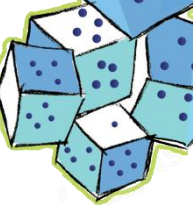# Generalized Tensor Decompositions for Non-Normal Data

## Tamara G. Kolda, Sandia Natl. Labs.
### Joint work with David Hong (Michigan), Jed Duersch (SNL)
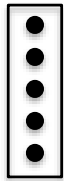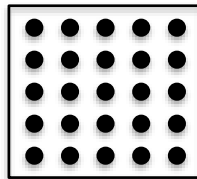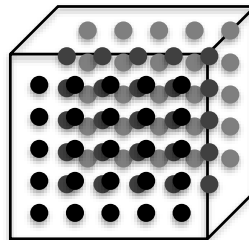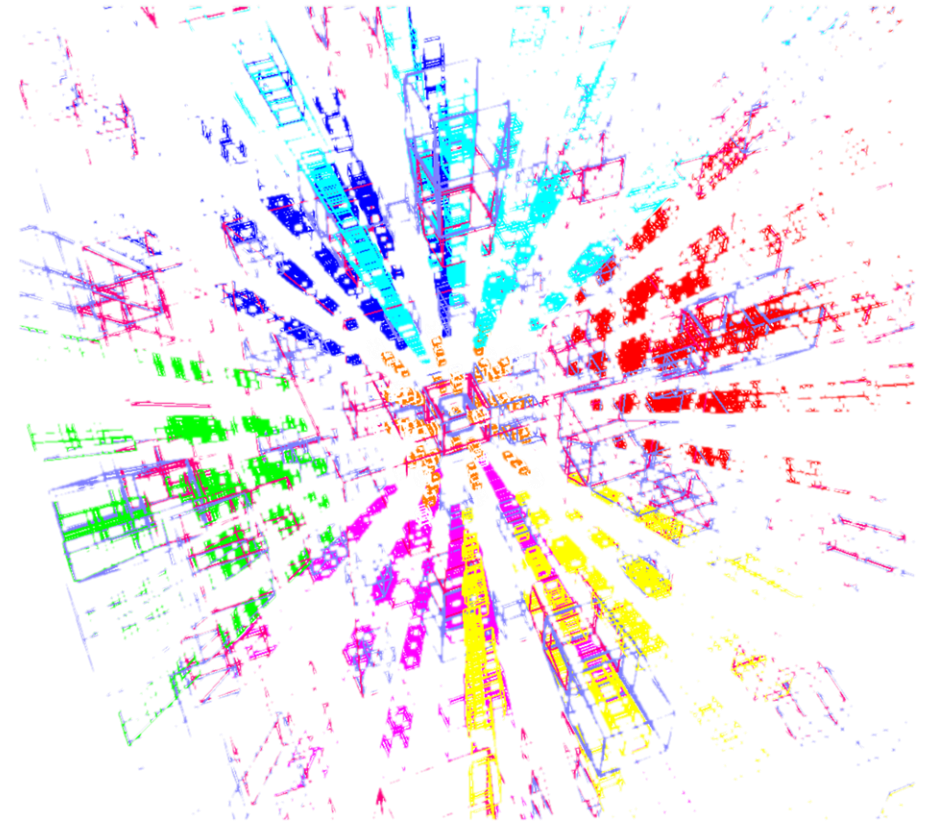
*Illustration by Chris Brigman*
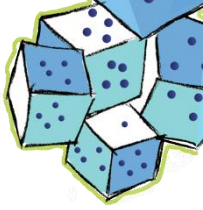
# A Tensor is an Multi-Way Array

Vector
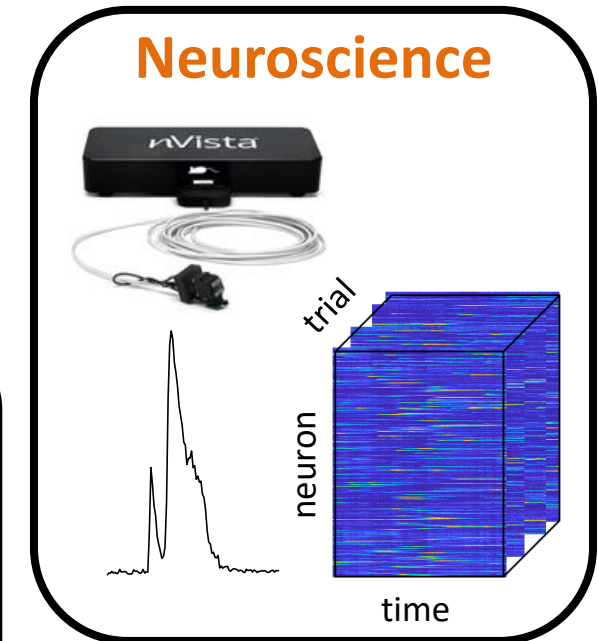$d = 1$

Matrix
$d = 2$

3rd-order Tensor
$d = 3$

$d^{th}$-order Tensor
$d > 3$

# Tensors Come From Many Applications

- **Chemometrics:** Emission x Excitation x Samples (Fluorescence Spectroscopy)

- **Neuroscience:** Neuron x Time x Trial (Calcium Imaging)

- **Criminology:** Day x Hour x Location x Crime (Chicago Crime Reports)

- **Medicine:** Channel x Wavelength x Time (EEG measurements)

- **Sports:** Player x Statistic x Season

- **Cyber-Traffic:** IP x IP x Port x Time

- **Social Network:** Person x Person x Time x Interaction-Type



**Chemometrics**

**Neuroscience**

**Criminology**

# Tensor Decomposition: A Mathematical & Statistical Tool for Analysis of Tensor Data

**Sandia National Laboratories**

Express the tensor as the sum of meaningful parts, which is the starting point for data analysis activities

**Data Analysis**

Includes visualization, clustering, filling in missing entries, etc.

**Sum of Parts**

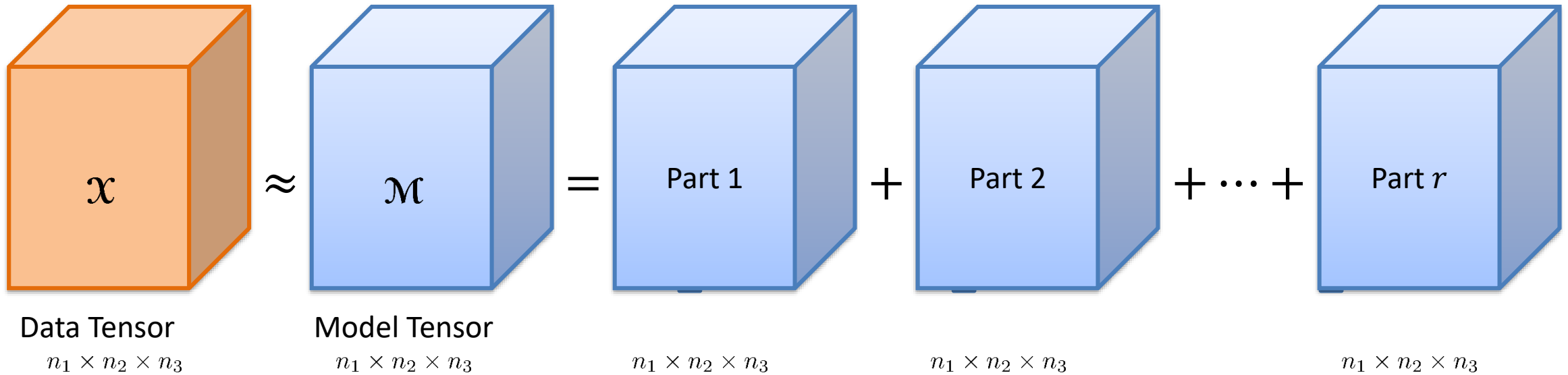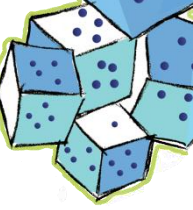**Mathematical & Statistical Tool**

Mathematics/Statistics play a role in....
- Defining the error metric
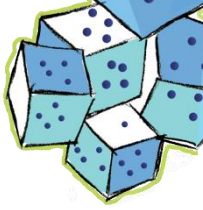- Developing efficient algorithms

## Related Concepts for Matrices

- Singular value decomposition (SVD)
- Principal component analysis (PCA)
- Independent component analysis (ICA)
- Nonnegative matrix factorization (NMF)
- Sparse matrix factorization
- Matrix completion

# Break Tensor into Understandable Parts…

$$\mathcal{X} \approx \mathcal{M} = \text{Part 1} + \text{Part 2} + \cdots + \text{Part } r$$

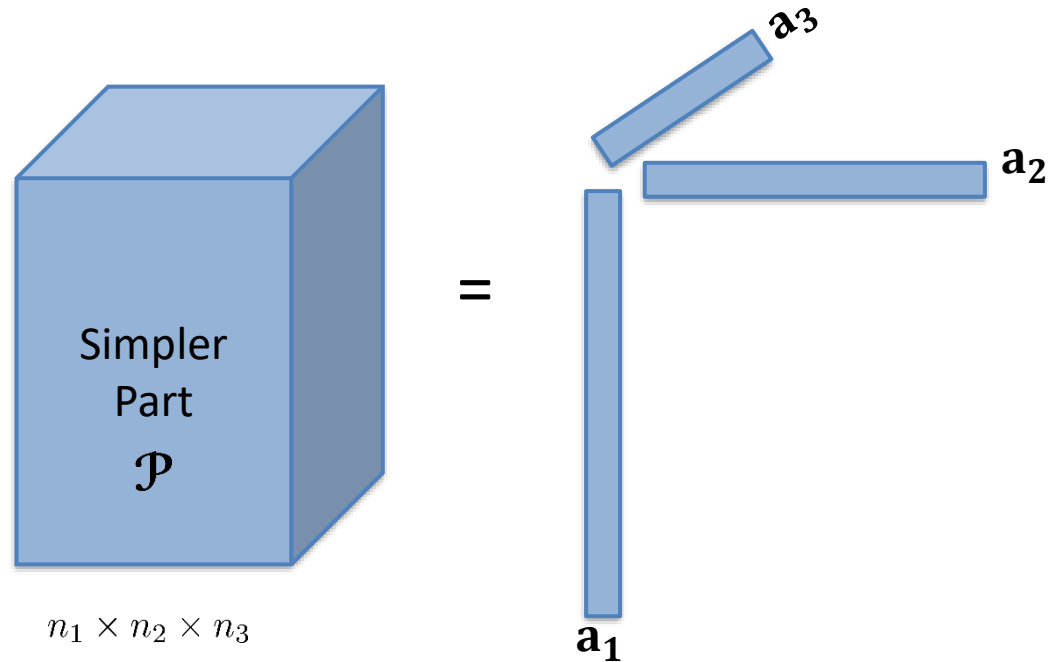| Data Tensor | Model Tensor | | | |
|---|---|---|---|---|
| $n_1 \times n_2 \times n_3$ | $n_1 \times n_2 \times n_3$ | $n_1 \times n_2 \times n_3$ | $n_1 \times n_2 \times n_3$ | $n_1 \times n_2 \times n_3$ |

Key: The parts have structure!

# Rank-1 Tensors are the "Parts"

Given **$d$ vectors**:

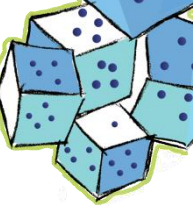$$\mathbf{a}_k \in \mathbb{R}^{n_k} \text{ for } k = 1, \ldots, d$$

The **outer product** is

$$\mathcal{P} = \mathbf{a}_1 \circ \mathbf{a}_2 \cdots \circ \mathbf{a}_d \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_d}$$
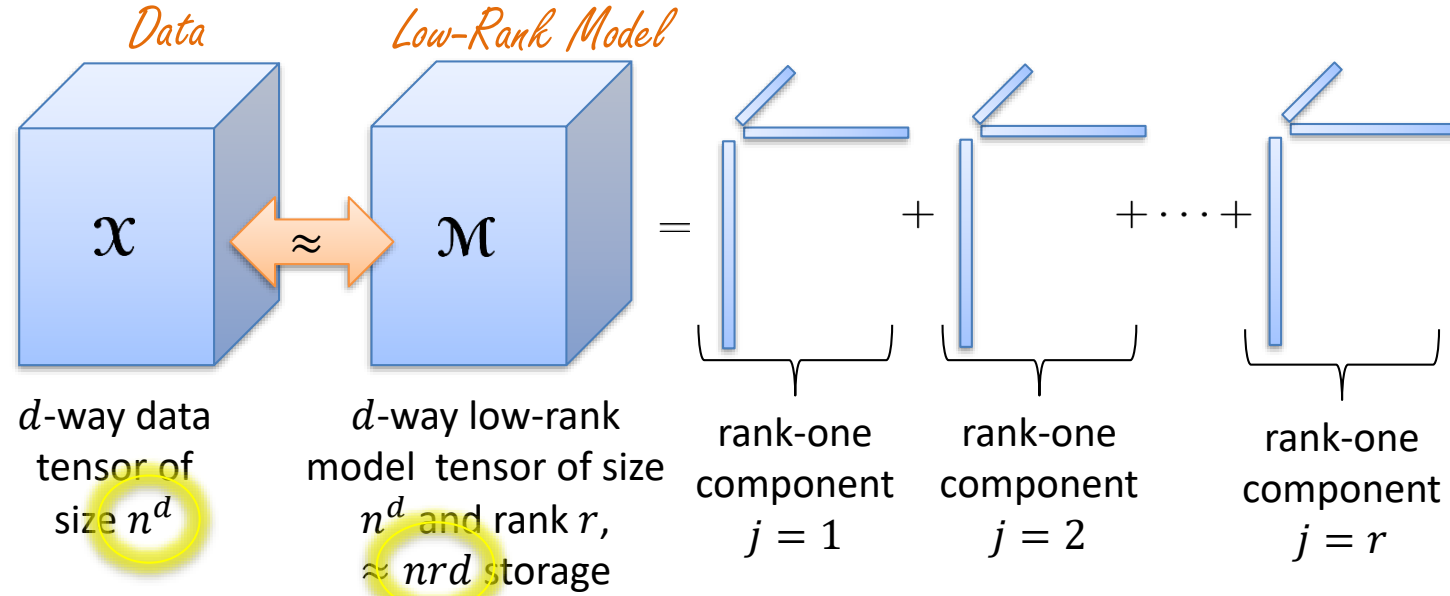


Simpler Part $\mathcal{P}$

$n_1 \times n_2 \times n_3$

$$\mathcal{P}(i_1, i_2, i_3) = \mathbf{a}_1(i_1)\,\mathbf{a}_2(i_2)\,\mathbf{a}_3(i_3)$$

# CANDECOMP/PARAFAC (CP) Tensor Factorization Uncovers the Rank-1 Parts

Images are three-way ($d = 3$), but assume all tensors are of size
$n_1 \times n_2 \times \cdots \times n_d$

WLOG, $n = n_1 = \cdots = n_d$

*Data*          *Low-Rank Model*

$$\mathcal{X} \approx \mathcal{M} = \sum_{j=1}^{r} + \cdots +$$

$d$-way data tensor of size $n^d$

$d$-way low-rank model tensor of size $n^d$ and rank $r$, $\approx nrd$ storage

rank-one component $j = 1$

rank-one component $j = 2$

rank-one component $j = r$

$$\mathcal{X} \approx \mathcal{M} \quad \text{where} \quad \mathcal{M} = \sum_{j=1}^{r} \mathbf{A}_1(:,j) \circ \mathbf{A}_2(:,j) \circ \cdots \circ \mathbf{A}_d(:,j)$$
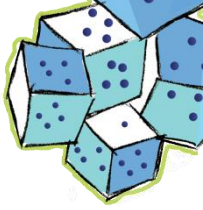
*Factor Matrices*

Low-rank: $\quad \text{rank}(\mathcal{M}) \leq r \ll n^d$

Factor matrices: $\quad \mathbf{A}_k \in \mathbb{R}^{n_k \times r} \text{ for } k \in \{1, \ldots, d\}$

Hitchcock, 1927; Carroll and Chang, 1970; Harshman, 1970

# CP first invented in 1927

Frank Lauren Hitchcock
MIT Professor
(1875–1957)



F. L. Hitchcock, *The Expression of a Tensor or a Polyadic as a Sum of Products*, Journal of Mathematics and Physics, 1927

# CP Independently Reinvented (twice) in 1970
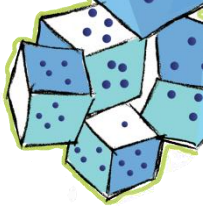


**CANDECOMP: <u>Can</u>onical <u>Decomp</u>osition**



J. Douglas Carroll
Bell Labs
(1939-2011)

Jih-Jie Chang
Bell Labs
(1927-2007)

**PARAFAC: <u>Para</u>llel <u>Fac</u>tors**



Richard A. Harshman
Univ. Ontario
(1943-2008)

**CP: CANDECOMP/PARAFAC**

In 2000, Henk Kiers proposed this *compromise* name

**CP: Canonical Polyadic**

2010: Pierre Comon, Lieven DeLathauwer, and others reverse-engineered CP, revising some of Hitchcock's terminology

*Many thanks to the following persons for helping me learn about Jih-Jie Chang: Fan Chung, Ron Graham, Shen Lin (husband), May Chang (niece), Lili Bruer (daughter).*

# Standard CP: Sum of Squares Error (SSE)

$$\mathcal{X} \approx \mathcal{M} = \left| \leftarrow + \right| \leftarrow + \cdots + \left| \leftarrow$$

Shorthand for element of data tensor:
$$x_i \equiv x(i_1, i_2, \ldots, i_d)$$

Element of model low-rank tensor:
$$m_i \equiv \sum_{j=1}^{r} \prod_{k=1}^{d} \mathbf{A}_k(i_k, j)$$

(defined in terms of factor matrices)

**Standard CP**

$$\min \ F(\mathcal{X}, \mathcal{M}) \equiv \sum_{i \in \Omega} (x_i - m_i)^2$$
$$\text{s.t.} \ \ \text{rank}(\mathcal{M}) \leq r$$

$\Omega$ = set of all $n^d$ elements in tensor

Hitchcock, 1927; Carroll and Chang, 1970; Harshman, 1970

# Generalized CP (GCP)



$$\mathcal{X} \approx \mathcal{M} = \left| \begin{array}{c} \\ \end{array} \right. + \left| \begin{array}{c} \\ \end{array} \right. + \cdots + \left| \begin{array}{c} \\ \end{array} \right.$$

**GCP**

$$\min \ F(\mathbf{\mathcal{X}}, \mathbf{\mathcal{M}}) \equiv \sum_{i \in \Omega} f(x_i, m_i)$$

$$\text{s.t. } \text{rank}(\mathbf{\mathcal{M}}) \leq r$$

## Why?

- SSE: maximum likelihood estimate (MLE) for Gaussian distribution

$$x_i = m_i + \epsilon, \ \epsilon \sim \mathcal{N}(0, \sigma)$$

$$x_i \sim \mathcal{N}(m_i, \sigma)$$

- Different MLEs for different distributions
  - Poisson (counts)
  - Bernoulli (binary)

Hong, Kolda, Duersch, SIAM Review, 2019

# Probability Distribution ⇒ Maximum Likelihood Estimator

Data Value

"Natural" Parameter

Model Value

$$x_i \sim p(x_i | \theta_i) \text{ where } \ell(\theta_i) = m_i$$

Link Function

Probability Distribution Function (PDF) or Probability Mass Function (PMF)

Maximize Likelihood of Data Tensor

$$\prod_{i \in \Omega} p(x_i, \theta_i)$$

Maximize Log-Likelihood

$$\sum_{i \in \Omega} \log p(x_i, \theta_i)$$

**GCP**

$$\min \; F(\mathcal{X}, \mathcal{M}) \equiv \sum_{i \in \Omega} f(x_i, m_i)$$

$$\text{s.t. } \operatorname{rank}(\mathcal{M}) \leq r$$

Given PDF/PMF $p(x|\theta)$ and link function $\ell(\theta)$, GCP MLE by minimizing

$$f(x, m) = -\log p(x, \ell^{-1}(m))$$

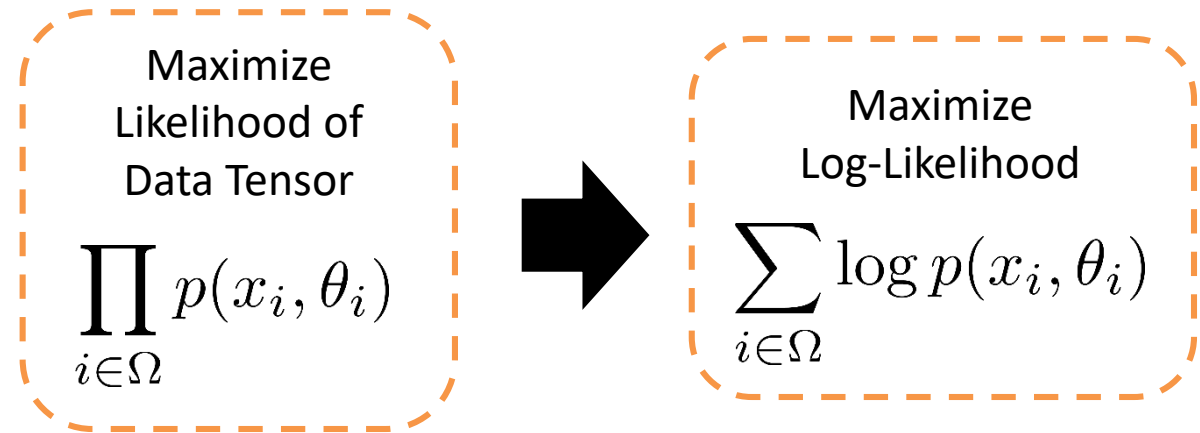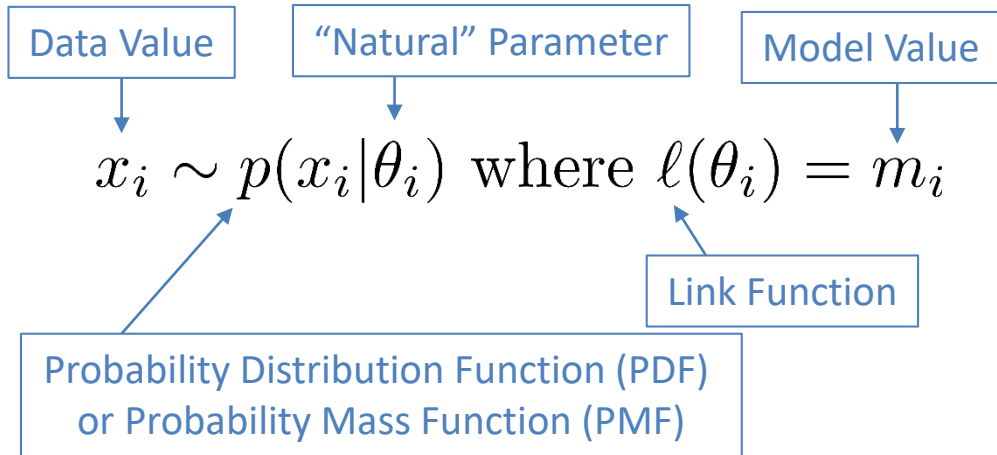Hong, Kolda, Duersch, SIAM Review, 2019

# Gaussian MLE (Standard CP)

**PDF for Normal Distribution**

$$p(x \mid \mu, \sigma) = \frac{e^{-(x-\mu)^2 / 2\sigma^2}}{\sqrt{2\pi\sigma^2}}$$

and

**Link Function**

$$m = \mu$$

$\sigma$ constant

Negative log-likelihood:

$$-\log p(x|\mu,\sigma) = \frac{(x-u)^2}{2\sigma^2} + \frac{1}{2}\log(2\pi\sigma^2)$$

Eliminate natural parameter via link function:

$$f(x,m) = \frac{(x-m)^2}{2\sigma^2} + \frac{1}{2}\log(2\pi\sigma^2)$$

Eliminate constants:

$$f(x,m) = (x-m)^2$$

Hong, Kolda, Duersch, SIAM Review, 2019

# Bernoulli MLE with Odds Link (Binary Data)

Bernoulli random variable

$$x \in \{0,1\}$$

$\rho$ = probability of a 1

$$p(x \mid \rho) = \rho^x (1 - \rho)^{(1-x)}, \quad x \in \{0, 1\}$$

**PMF for Bernoulli Distribution**

$$p(x \mid \rho) = \rho^x (1 - \rho)^{(1-x)}$$
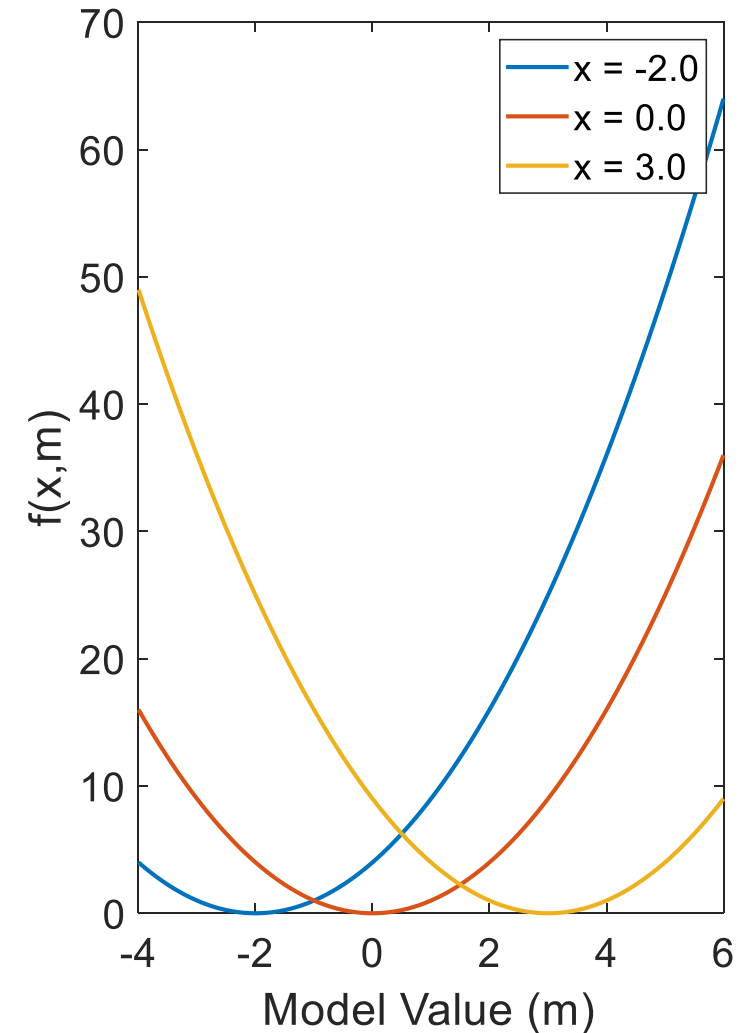$$x \in \{0, 1\}$$

and

**Link Function**

$$m = \frac{\rho}{(1 - \rho)}$$

**Odds Link**

$$\ell(\rho) = \rho / (1 - \rho)$$
$$\ell^{-1}(m) = m / (1 + m)$$

| Odds ($m$) | Probability ($\rho$) |
|:---:|:---:|
| ¼ | 20% |
| 1 | 50% |
| 4 | 80% |
| 10 | 90% |

Negative log-likelihood:

$$-\log p(x \mid \rho) = \log \frac{1}{1 - \rho} - x \log \frac{\rho}{1 - \rho}$$

Eliminate natural parameter
via link function:

$$f(x, m) = \log(1 + m) - x \log m \quad \text{for} \quad m > 0$$

Hong, Kolda, Duersch, SIAM Review, 2019

# Bernoulli MLE with Odds Link (Binary Data)

Bernoulli random variable

$$x \in \{0,1\}$$

$$\rho = \text{probability of a 1}$$

$$p(x \mid \rho) = \rho^x (1 - \rho)^{(1-x)}, \quad x \in \{0, 1\}$$

PMF for Bernoulli Distribution

$$p(x \mid \rho) = \rho^x (1 - \rho)^{(1-x)}$$

$$x \in \{0, 1\}$$
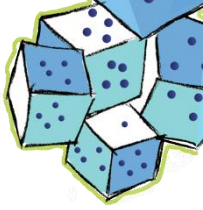
and

Link Function

$$m = \frac{\rho}{(1 - \rho)}$$

Negative log-likelihood:

$$-\log p(x \mid \rho) = \log \frac{1}{1 - \rho} - x \log \frac{\rho}{1 - \rho}$$

Eliminate natural parameter
via link function:

$$f(x, m) = \log(1 + m) - x \log m \quad \text{for} \quad m > 0$$

Hong, Kolda, Duersch, SIAM Review, 2019

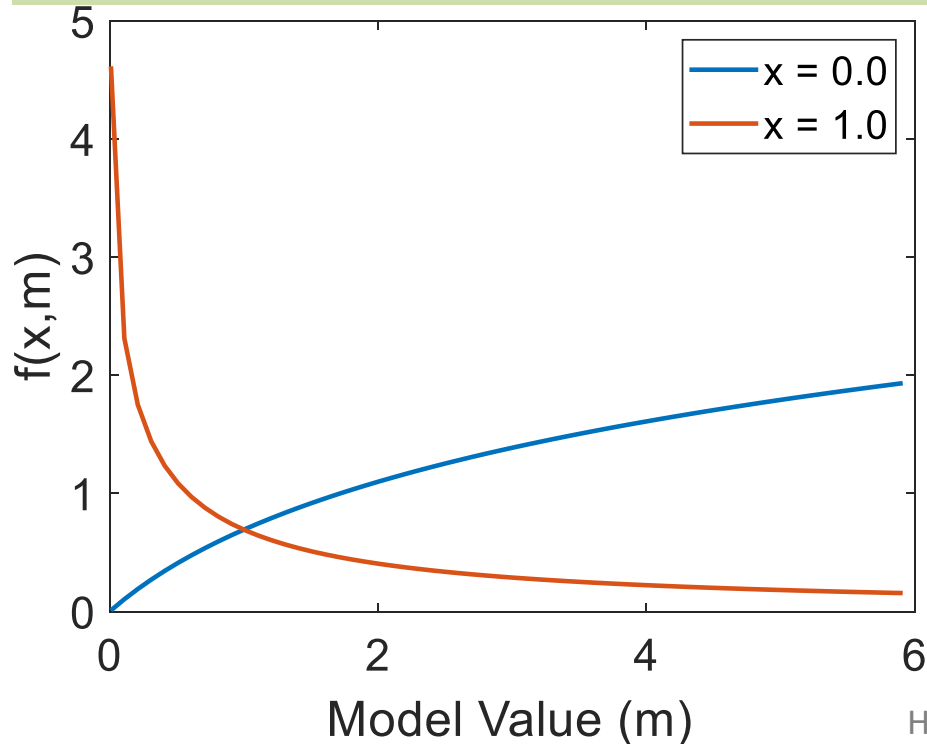# Bernoulli MLE with Logit Link (Binary Data)

Bernoulli random variable
$$x \in \{0,1\}$$
$$\rho = \text{probability of a 1}$$

$$p(x \mid \rho) = \rho^x (1-\rho)^{(1-x)}, \quad x \in \{0,1\}$$

**PMF for Bernoulli Distribution**
$$p(x \mid \rho) = \rho^x (1-\rho)^{(1-x)}$$
$$x \in \{0,1\}$$

and

**Link Function**
$$m = \log \frac{\rho}{(1-\rho)}$$

**Logit (Log-Odds) Link**
$$\ell(\rho) = \log\big(\rho / (1-\rho)\big)$$
$$\ell^{-1}(m) = e^m / (1+e^m)$$

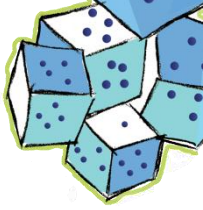| Log-Odds($m$) | Probability ($\rho$) |
|---|---|
| -1.39 | 20% |
| 0 | 50% |
| 1.39 | 80% |
| 2.30 | 90% |

Negative log-likelihood:

$$-\log p(x \mid \rho) = \log \frac{1}{1-\rho} - x \log \frac{\rho}{1-\rho}$$

Eliminate natural parameter
via link function:

$$f(x, m) = \log(1 + e^m) - xm \quad \text{for} \quad m \in \mathbb{R}$$

Hong, Kolda, Duersch, SIAM Review, 2019
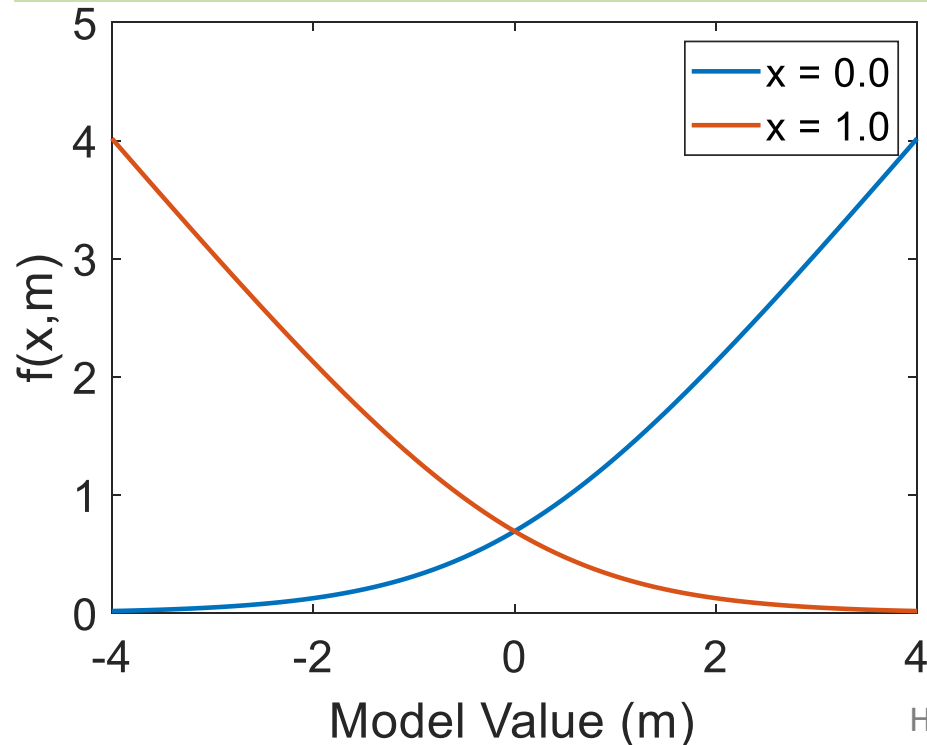
# Bernoulli MLE with Logit Link (Binary Data)

**Bernoulli random variable**
$$x \in \{0,1\}$$
$$\rho = \text{probability of a 1}$$
$$p(x \mid \rho) = \rho^x (1-\rho)^{(1-x)}, \quad x \in \{0,1\}$$

**PMF for Bernoulli Distribution**
$$p(x \mid \rho) = \rho^x (1-\rho)^{(1-x)}$$
$$x \in \{0,1\}$$

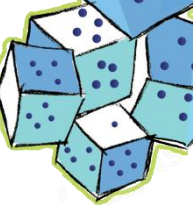and

**Link Function**
$$m = \log \frac{\rho}{(1-\rho)}$$

Negative log-likelihood:

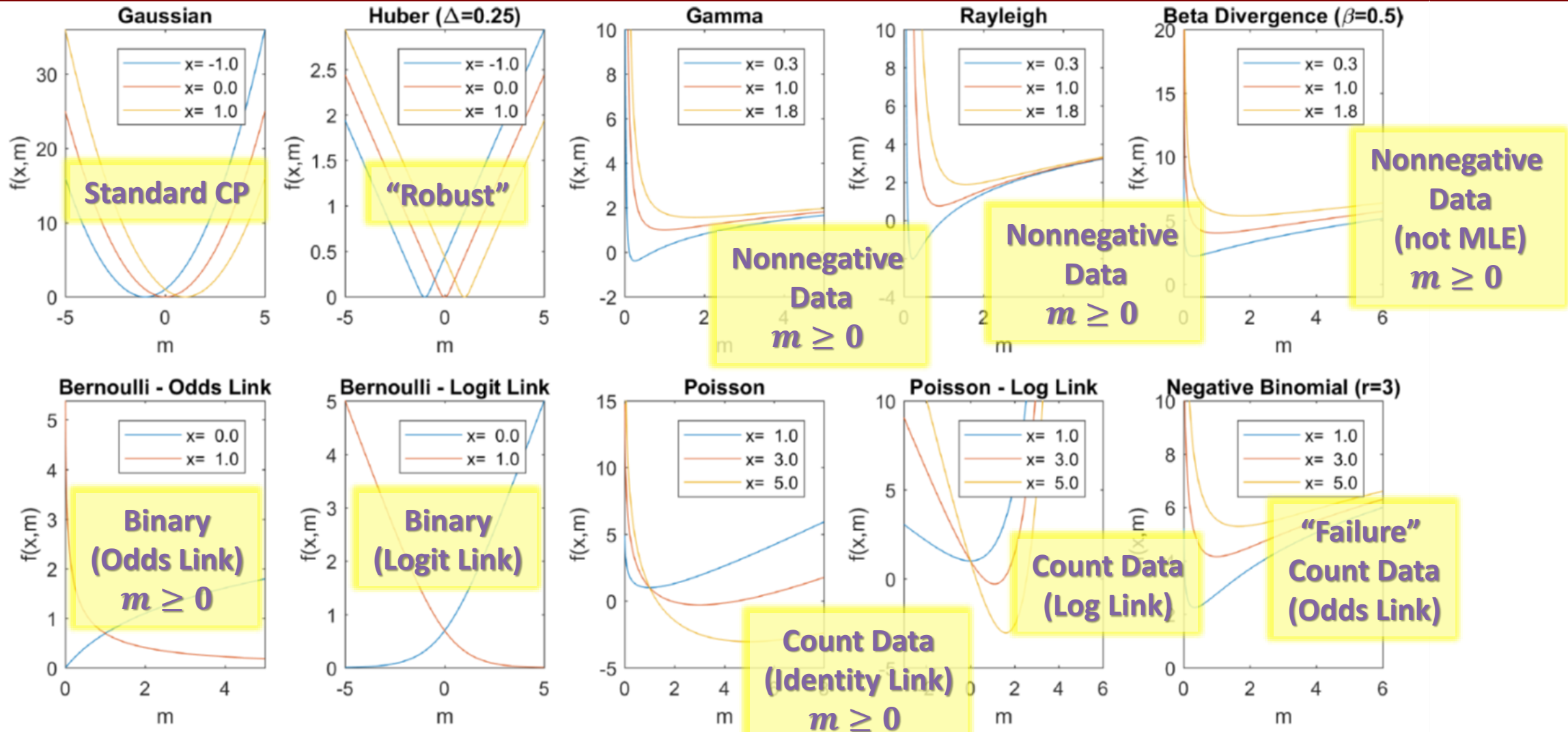$$-\log p(x \mid \rho) = \log \frac{1}{1-\rho} - x \log \frac{\rho}{1-\rho}$$

Eliminate natural parameter
via link function:

$$f(x,m) = \log(1 + e^m) - xm \quad \text{for} \quad m \in \mathbb{R}$$

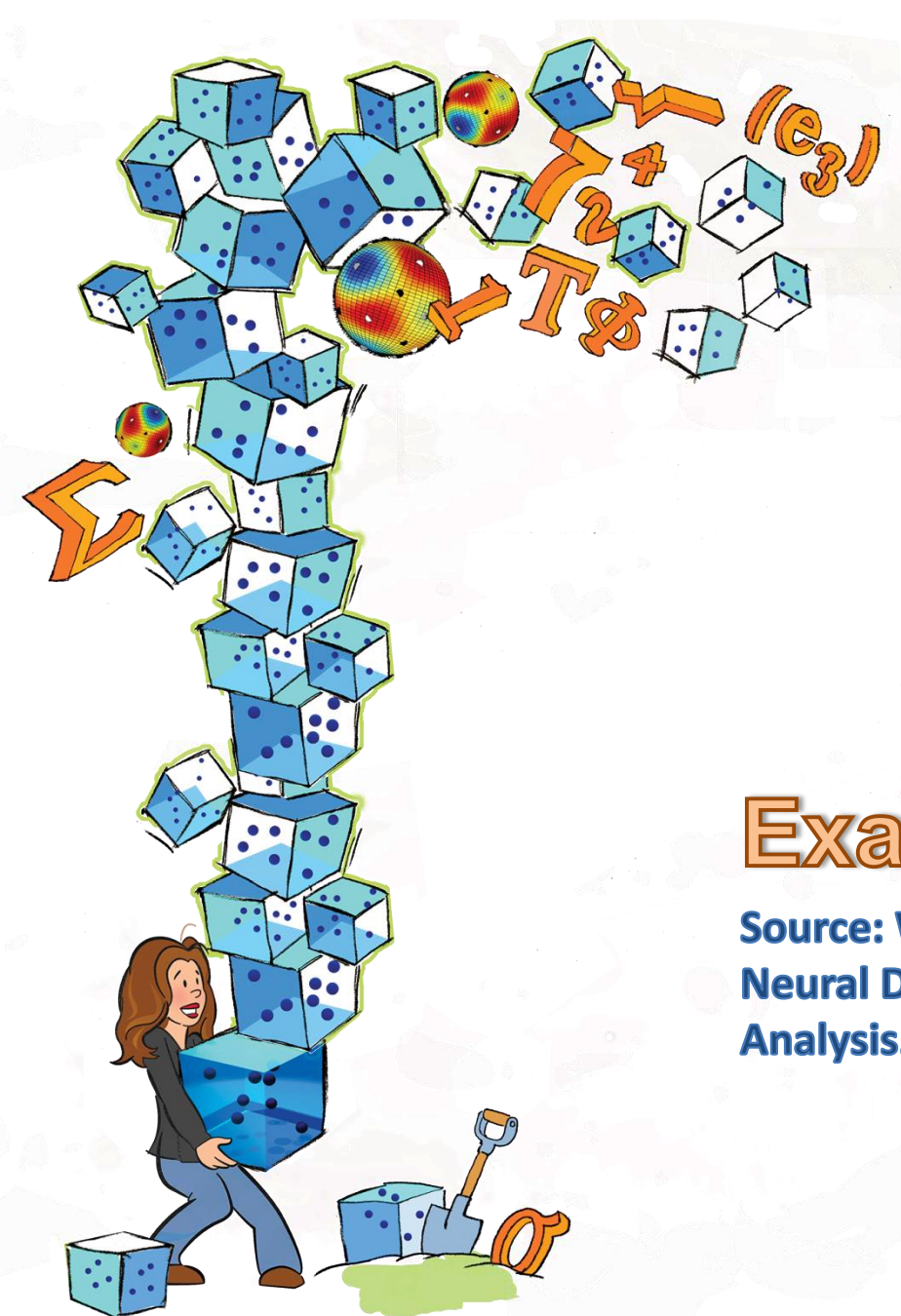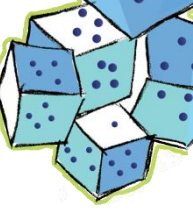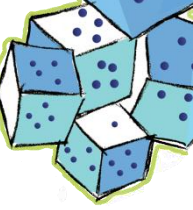f(x,m) vs Model Value (m), with legend: x = 0.0, x = 1.0

Hong, Kolda, Duersch, SIAM Review, 2019

# Sampling of Loss Functions



Hong, Kolda, Duersch, SIAM Review, 2019

# Example Tensor from Neuroscience

# Activity of Single Neuron Measured Over Time Produces Vector Data

*Thanks to Schnitzer Group @ Stanford*
Mark Schnitzer, Fori Wang, Tony Kim

111 time bins

Microscope by Inscopix



mouse in maze

neural activity via calcium imaging

Williams et al., Neuron, 2018

# Multiple Neurons Measured Over Time Produces Matrix

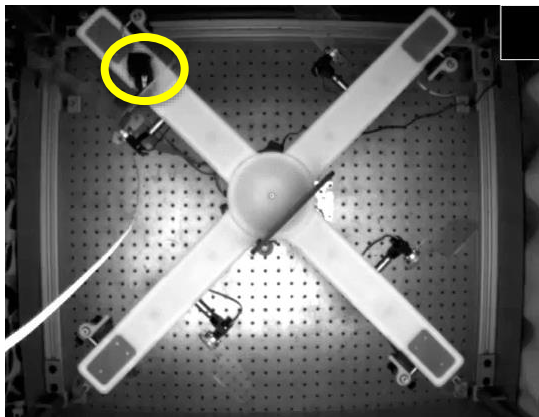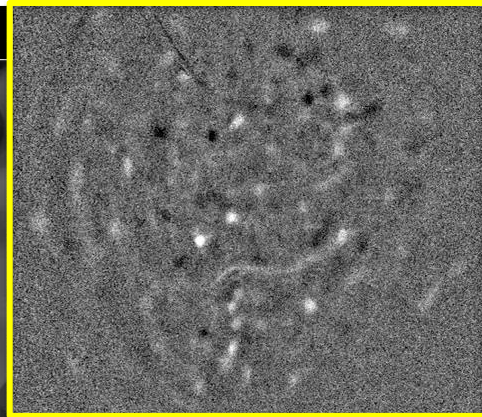*Thanks to Schnitzer Group @ Stanford*
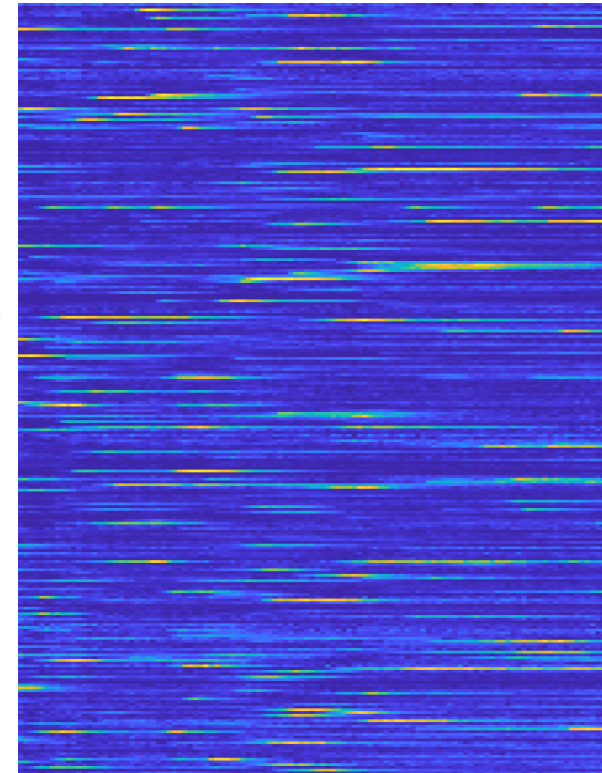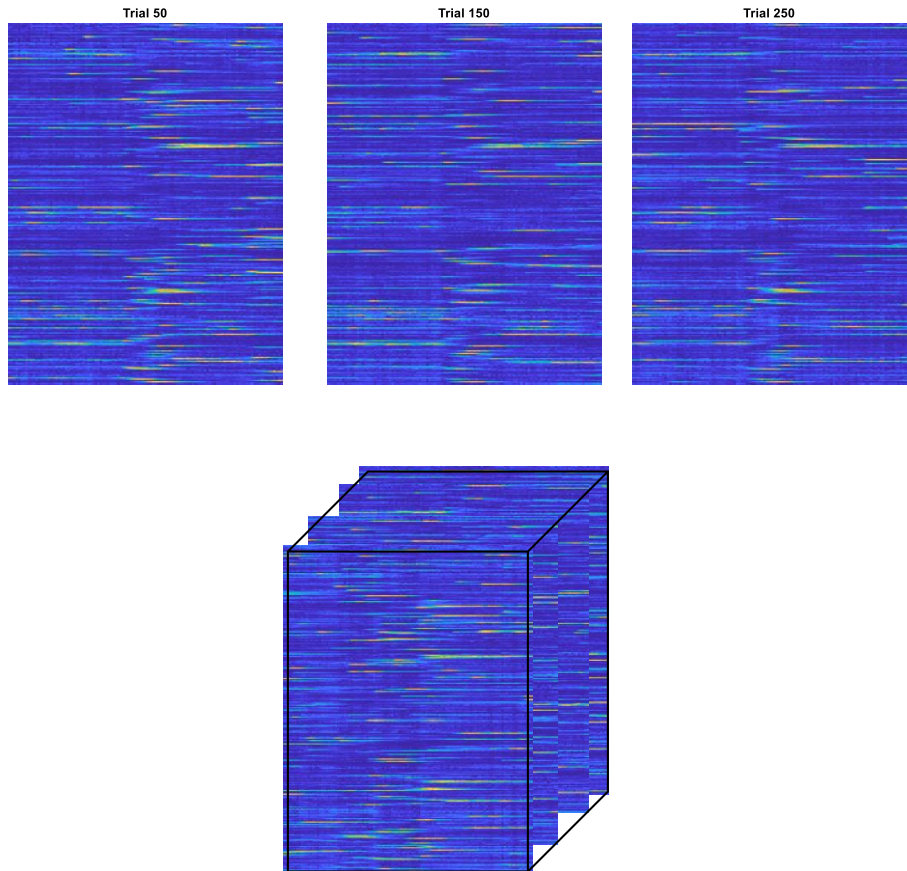Mark Schnitzer, Fori Wang, Tony Kim

Microscope by
Inscopix

282 neurons × 111 time bins

mouse
in "maze"

neural activity



Williams et al., Neuron, 2018

# Multiple Trials Produces 3-way Tensor

Trial 50    Trial 150    Trial 250

282 neurons × 111 time bins × 300 trials



- 300 Trials over 5 Days
- Start West
- Conditions Swap Twice
  - Turn South
  - Turn North
  - Turn South

Williams et al., Neuron, 2018

# Example Neuron Activity

Thin lines show 300 individual trials

Thick line is average



Hong, Kolda, Duersch, SIAM Review, 2019

Neuron Modes Plotted as a Bar Chart
(Red Lines Correspond to Examples in Previous Slide)

Hong, Kolda, Duersch, SIAM Review, 2019

# Time Factor Vector Visualized as Line



Time (within trial) Plotted as a Line
(Dashed Line is Zero)

Hong, Kolda, Duersch, SIAM Review, 2019

Rule Change

Rule Change

Trial Plotted as Scatter Graph
Right turn = Green
Left turn = Orange
Filled = Reward

Hong, Kolda, Duersch, SIAM Review, 2019

Hong, Kolda, Duersch, SIAM Review, 2019

# "Standard" CP Decomposition of Mouse Data, aka Gaussian ($f(x, m) = (x - m)^2$)

Neuron (scaled) | Time | Trial (Green/Orange = Turn Right/Left, Reward = Filled)

Turn Direction

# CP Tensor Decomposition Can be Tough to Interpret due to Negative Entries

# Regression Using GCP Factors on Trial Mode

Trial Factor Matrix is $300 \times 8$



$$\min_{\boldsymbol{\beta}} \| \mathbf{A}_3^{\mathrm{train}} \boldsymbol{\beta} - \mathbf{y}^{\mathrm{train}} \|$$

$$\hat{\mathbf{y}}^{\mathrm{test}} = \left[ \mathbf{A}_3^{\mathrm{test}} \boldsymbol{\beta} \geq 0.5 \right]$$

Look at predicting turn and reward.
Split into two groups of 150 trials.
Train regression model with 1st group.
Test with 2nd group.
Repeat 100 times.

Regression Errors in 100
Trials (15000 predictions)



Turn  Reward

Hong, Kolda, Duersch, SIAM Review, 2019

# Optimization Formulation for GCP Tensor Decomposition

**GCP**

$$\min \ F(\mathbf{X}, \mathbf{M}) \equiv \sum_{i \in \Omega} f(x_i, m_i)$$

$$\text{s.t.} \ \text{rank}(\mathbf{M}) \leq r$$

$i$ = multi-index
$\Omega$ = all indices

- Standard CP [Hitchcock, 1927; Carrol & Chang, 1970; Harshman, 1970]

$$f(x, m) = (x - m)^2$$

- Poisson CP (Identity Link) [Welling & Webber, 2001; Chi & Kolda, 2009]

$$f(x, m) = m - x \log m$$

- Logistic CP, etc. [Hong, Kolda, Duersch, 2018]

$$f(x, m) = \log(m + 1) - x \log(m)$$



$\mathbf{X} \approx \mathbf{M}$

$d$-way data tensor of size $n^d$

$d$-way low-rank model tensor of size $n^d$ and rank $r$

rank-one component $j = 1$

rank-one component $j = 2$

rank-one component $j = r$

$$\mathbf{X} \approx \mathbf{M} \quad \text{where} \quad \mathbf{M} = \sum_{j=1}^{r} \mathbf{A}_1(:, j) \circ \mathbf{A}_2(:, j) \circ \cdots \circ \mathbf{A}_d(:, j)$$

Low-rank: $\quad \text{rank}(\mathbf{M}) \leq r \ll n^d$

Factor matrices: $\quad \mathbf{A}_k \in \mathbb{R}^{n_k \times r} \text{ for } k \in \{1, \ldots, d\}$

WLOG, $n = n_1 = \cdots = n_d$

# Gradient-based Optimization for Fitting the GCP Model

**GCP**

$$\min \; F(\mathbf{\mathcal{X}}, \mathbf{\mathcal{M}}) \equiv \sum_{i \in \Omega} f(x_i, m_i)$$

$$\text{s.t. } \operatorname{rank}(\mathbf{\mathcal{M}}) \leq r$$

Gradients computed via a sequence of matricized-tensor times Khatri-Rao product (MTTKRPs):

$$\mathbf{G}_k \equiv \frac{\partial F}{\partial \mathbf{A}_k} = \mathbf{Y}_{(k)} \mathbf{Z}_k \text{ for } k = 1, \ldots, d$$

**MTTKRP**

gradient for mode $k$ factor matrix of size $n \times r$

tensor unfolded in mode $k$ into matrix of size $n \times n^{d-1}$

<u>Define</u>: Elementwise partial gradient tensor, same size as data tensor = $n^d$

$\mathbf{\mathcal{Y}}$ $\qquad y_i = \frac{\partial f}{\partial m}(x_i, m_i)$

<u>Define</u>: Khatri-Rao product in all modes but one of size $n^{d-1} \times r$

$$\mathbf{Z}_k = \mathbf{A}_d \odot \cdots \odot \mathbf{A}_{k+1} \odot \mathbf{A}_{k-1} \odot \cdots \odot \mathbf{A}_1$$

MTTKRPs can be computed efficiently…
- Bader & Kolda, SISC, 2007 – Dense and sparse
- Phan, Tichavsky, Cichocki, 2013 – Sequence
- Smith et al., IPDPS 2015 – Sparse
- Kaya & Ucar, SC 2015 – Sparse
- Li et al., IPDPS 2017 – Sparse
- Hayashi et al., 2017 – Dense
- Ballard, Knight, Rouse, 2017 – Dense

# Stochastic Gradient Descent (SGD) for GCP

**30-Second Tutorial on SGD**

$$\min F(x)$$

Gradient Descent (GD)
$\alpha$ = learning rate
$$x^{(t+1)} = x^{(t)} - \alpha g^{(t)}$$

Stochastic Gradient Descent (SGD)
$$x^{(t+1)} = x^{(t)} - \alpha \tilde{g}^{(t)}$$
$$\mathbb{E}[\tilde{g}^{(t)}] = g^{(t)} \equiv \nabla F(x^{(t)})$$

Adam (Kingma & Ba, 2015)
*Adaptive momentum SGD*

---

**Standard gradient**     $\mathbf{G}_k = \mathbf{Y}_{(k)}\mathbf{Z}_k$     Cost: $O(rn^d)$ flops

$\mathcal{Y}$

$$y_i = \frac{\partial f}{\partial m}(x_i, m_i)$$

---

**Stochastic gradient**     $\tilde{\mathbf{G}}_k = \tilde{\mathbf{Y}}_{(k)}\mathbf{Z}_k$     Cost: $O(rs)$ flops

$\tilde{\mathcal{Y}}$

Choose stochastic *sparse* Y-tensor

$$\mathbb{E}[\tilde{\mathcal{Y}}] = \mathcal{Y}$$

such that

$$\mathrm{nnz}(\tilde{\mathcal{Y}}) \leq s \ll n^d$$

By linearity of expectation: $\mathbb{E}[\tilde{\mathbf{G}}_k] = \mathbf{G}_k$

# Uniform Sampling

## Goal: Random *sparse* tensor of size $n^d$ that equals the "Y-tensor" in expectation

Sample $s \ll n^d$ random tensor entries (with replacement)

$\tilde{s}_i = \#$ times $i$ sampled

$\tilde{y}_i = \tilde{s}_i \cdot \dfrac{n^d}{s} \cdot y_i$

$$\sum_{i \in \Omega} \tilde{s}_i = s$$

$y_i = \dfrac{\partial f}{\partial m}(x_i, m_i)$

**Theory**

Claim: $\mathbb{E}[\tilde{\mathcal{Y}}] = \mathcal{Y}$

Proof: $\mathbb{E}[\tilde{s}_i] = \dfrac{s}{n^d}$

$\mathbb{E}[\tilde{y}_i] = \mathbb{E}[\tilde{s}_i] \cdot \dfrac{n^d}{s} \cdot y_i = y_i$

Choosing $s$, the number of sampled elements…

* Choose $s = O(n)$

* Gradient = $O(rs) = O(rn)$ versus $O(rn^d)$

Downside…

* If data tensor is sparse, few entries corresponding to nonzeros will be chosen

# Stratified 0/1 Sampling
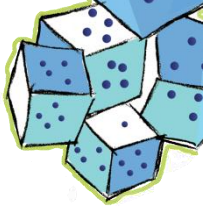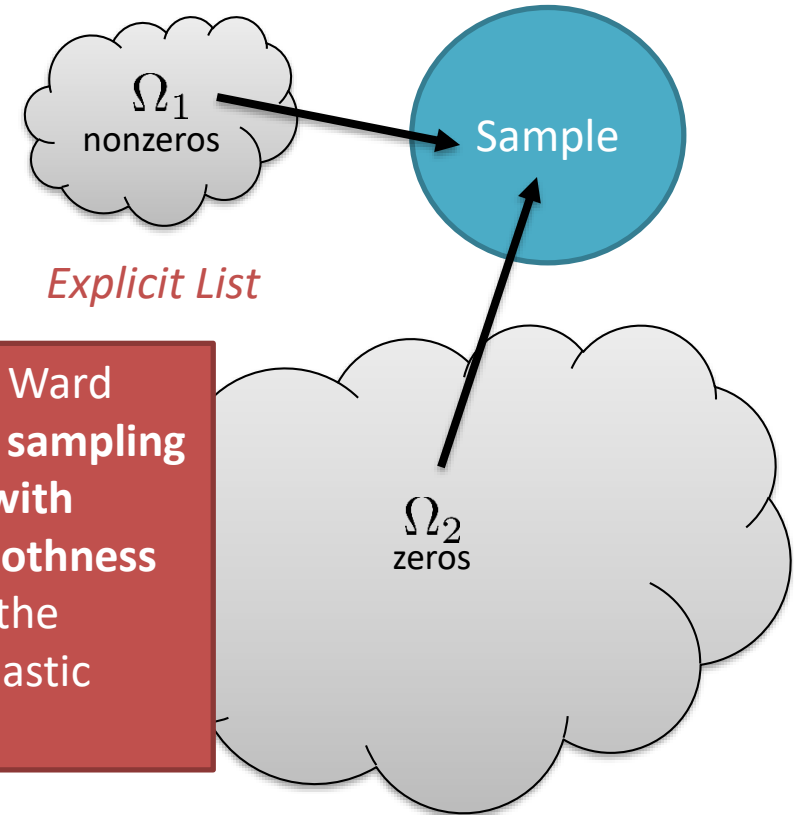
**Goal**: Random *sparse* tensor of size $n^d$ that equals the "Y-tensor" in expectation

Sample $p$ nonzeros and $q$ zeros.

$$\tilde{p}_i = \# \text{ times nonzero } i \text{ sampled} \qquad \eta = \# \text{ nonzeros}$$
$$\tilde{q}_i = \# \text{ times zero } i \text{ sampled} \qquad \zeta = \# \text{ zeros}$$
$$\tilde{y}_i = \left( \tilde{p}_i \cdot \frac{\eta}{p} + \tilde{q}_i \cdot \frac{\zeta}{q} \right) \cdot y_i$$

$$y_i = \frac{\partial f}{\partial m}(x_i, m_i)$$

**Theory**

Claim: $\mathbb{E}[\tilde{\mathcal{Y}}] = \mathcal{Y}$

Proof: $\mathbb{E}[\tilde{p}_i] = \dfrac{p}{\eta}, \ \mathbb{E}[\tilde{q}_i] = \dfrac{q}{\zeta}$

$$x_i = 1 \Rightarrow \mathbb{E}[\tilde{y}_i] = \mathbb{E}[\tilde{p}_i] \cdot \frac{\eta}{p} \cdot y_i = y_i$$

$$x_i = 0 \Rightarrow \mathbb{E}[\tilde{y}_i] = \mathbb{E}[\tilde{q}_i] \cdot \frac{\zeta}{q} \cdot y_i = y_i$$

$\Omega_1$ nonzeros

Sample

*Explicit List*

Needell, Srebro, and Ward (2013) justify **biased sampling toward functionals with higher Lipschitz smoothness** constants to reduce the variance in the stochastic gradient.
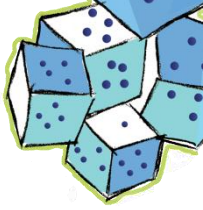
$\Omega_2$ zeros

*Implicit List (Requires Rejection Sampling)*

# Semi-Stratified 0/1 Sampling
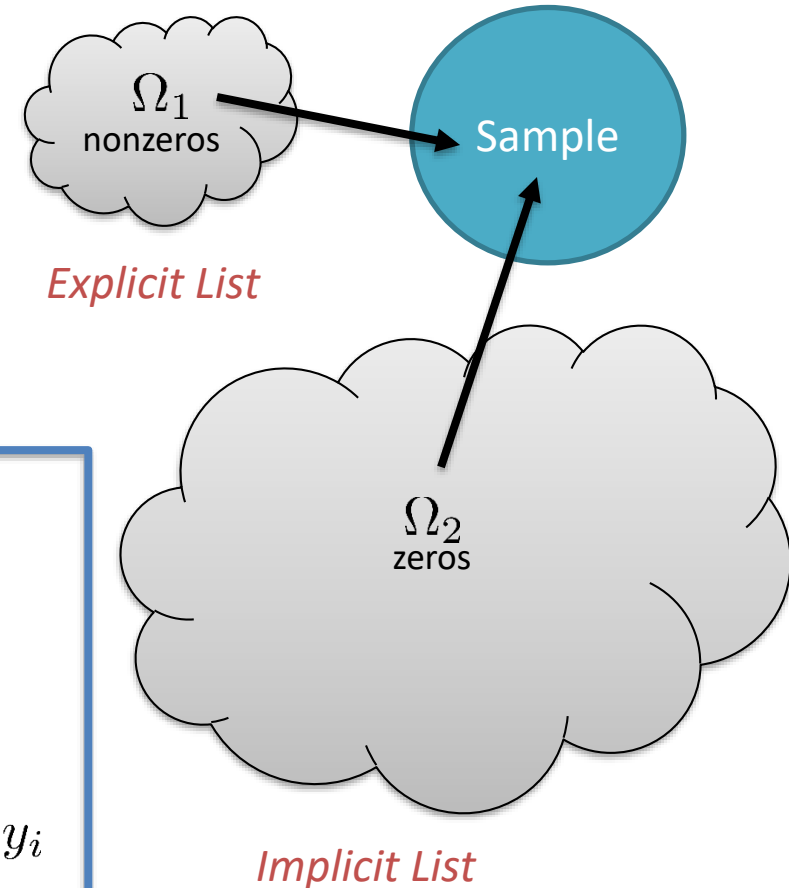
**Sandia National Laboratories**

**Goal**: Random *sparse* tensor of size $n^d$ that equals the "Y-tensor" in expectation

Sample $p$ nonzeros and $q$ **assumed** zeros.

$$\tilde{p}_i = \# \text{ times nonzero } i \text{ sampled} \qquad \eta = \# \text{ nonzeros}$$
$$\tilde{q}_i = \# \text{ times "zero" } i \text{ sampled} \qquad \zeta = \# \text{ zeros}$$

$$\tilde{y}_i = \tilde{p}_i \cdot \frac{\eta}{p} \cdot (y_i - c_i) + \tilde{q}_i \cdot \frac{(\eta + \zeta)}{q} \cdot c_i \text{ with } c_i \equiv \frac{\partial f}{\partial m}(0, m_i)$$

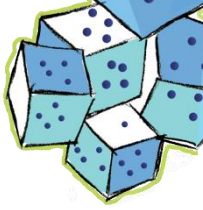$$y_i = \frac{\partial f}{\partial m}(x_i, m_i)$$

$\Omega_1$ nonzeros

*Explicit List*

Sample

**Theory**

Claim: $\mathbb{E}[\tilde{\mathcal{Y}}] = \mathcal{Y}$

Proof: $\mathbb{E}[\tilde{p}_i] = \dfrac{p}{\eta}, \ \mathbb{E}[\tilde{q}_i] = \dfrac{q}{(\zeta + \eta)}$
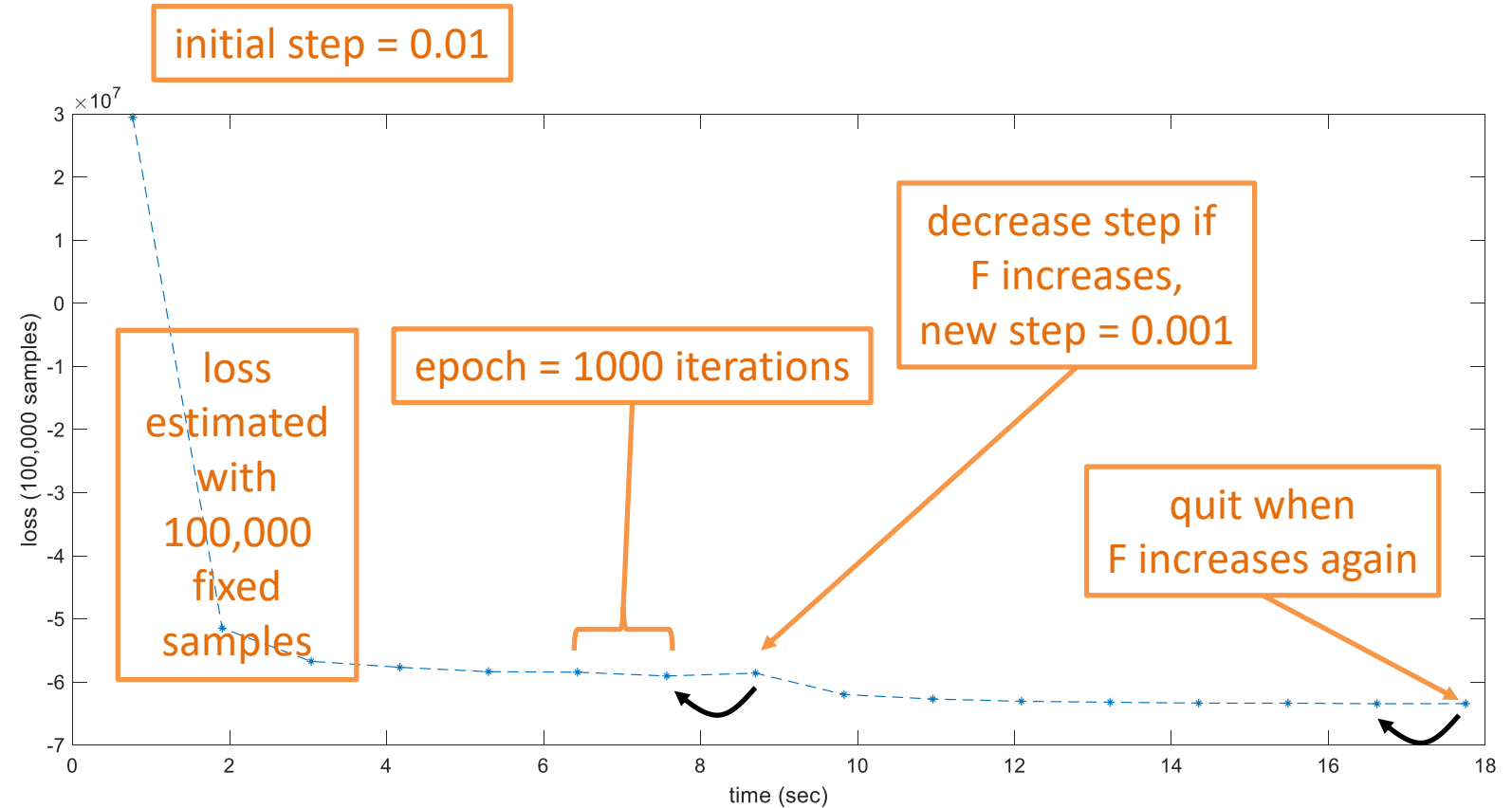
$$x_i = 0 \Rightarrow \mathbb{E}[\tilde{y}_i] = \mathbb{E}[\tilde{q}_i] \cdot \frac{(\eta + \zeta)}{q} \cdot y_i = y_i$$
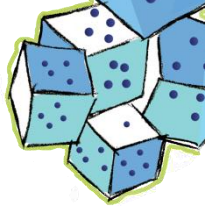
$$x_i = 1 \Rightarrow \mathbb{E}[\tilde{y}_i] = \mathbb{E}[\tilde{p}_i] \cdot \frac{\eta}{p} \cdot (y_i - c_i) + \mathbb{E}[\tilde{q}_i] \cdot \frac{\eta + \zeta}{q} \cdot c_i = y_i$$

$\Omega_2$ zeros

*Implicit List*
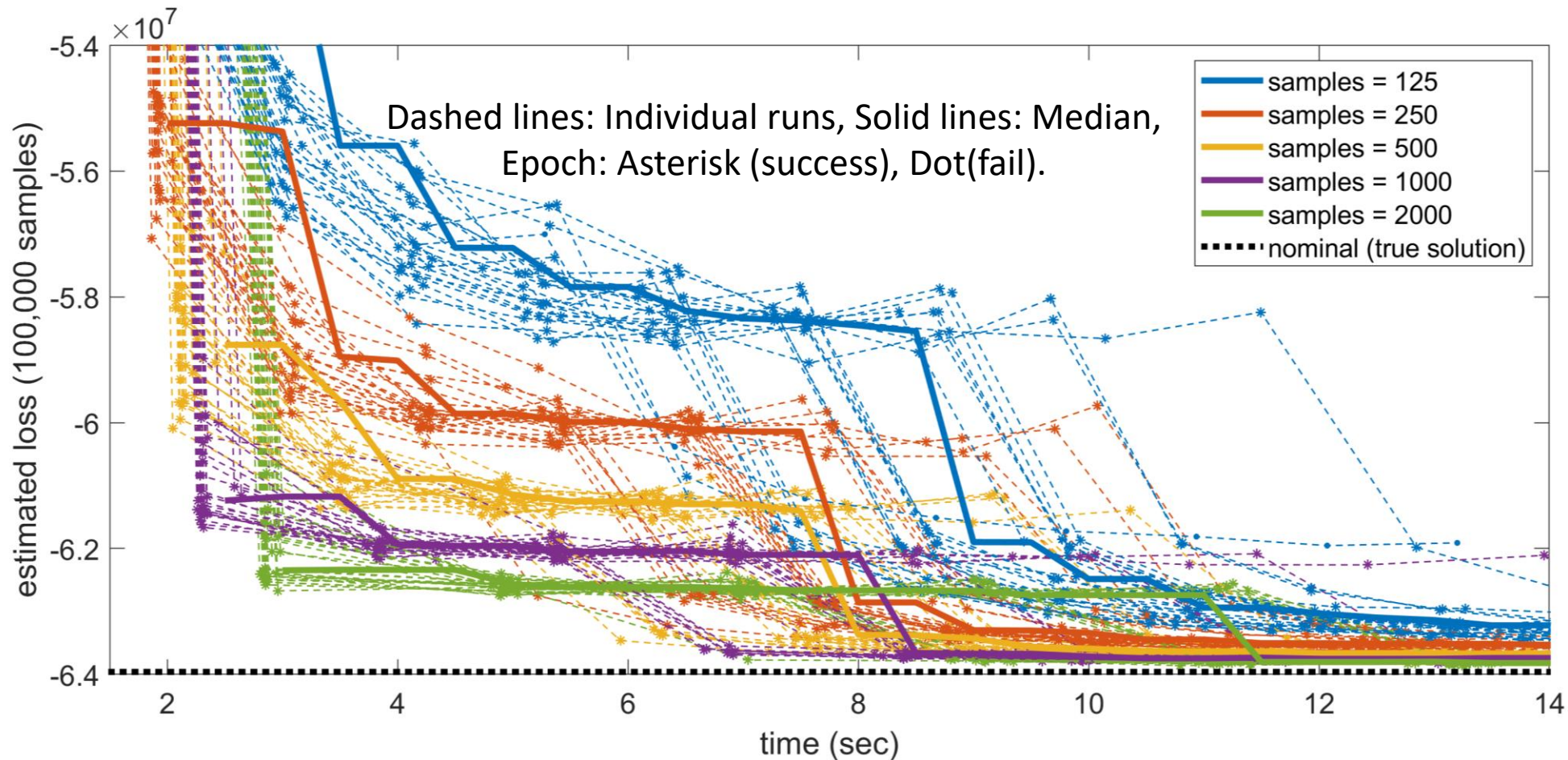
# GCP with Stochastic Optimization

- Nonconvex problem
  - No guarantees of finding minimizer
- Using Adam (Kingma & Ba, 2015)
  - Default parameters
  - Some tweaks for checking convergence
- Past work on recommender systems uses SGD but ignores zeros
  - Gemulla, Nijkamp, Hass, Sismanis, KDD'11
  - Zhuang, Chin, Juan, and Lin, RecSys'13
- Past work on streaming uses SGD but data appears one slice at a time
  - Mardani, Mateos, Giannakis, IEEE TSP 2015
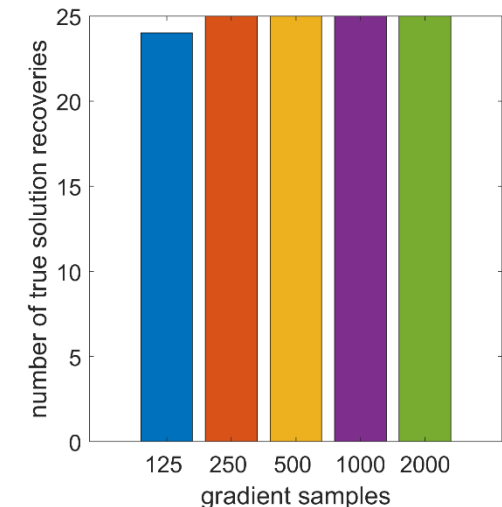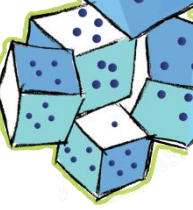  - Maehara, Hayashi, Kawarabayashi,

initial step = 0.01

decrease step if
F increases,
new step = 0.001

loss estimated with 100,000 fixed samples

epoch = 1000 iterations

quit when F increases again

# Example on Gamma-Distributed Data

$200 \times 150 \times 100 \times 50$ Tensor with low-rank ($r = 5$) structure based on Gamma distribution ($k = 1, \theta$ from model).

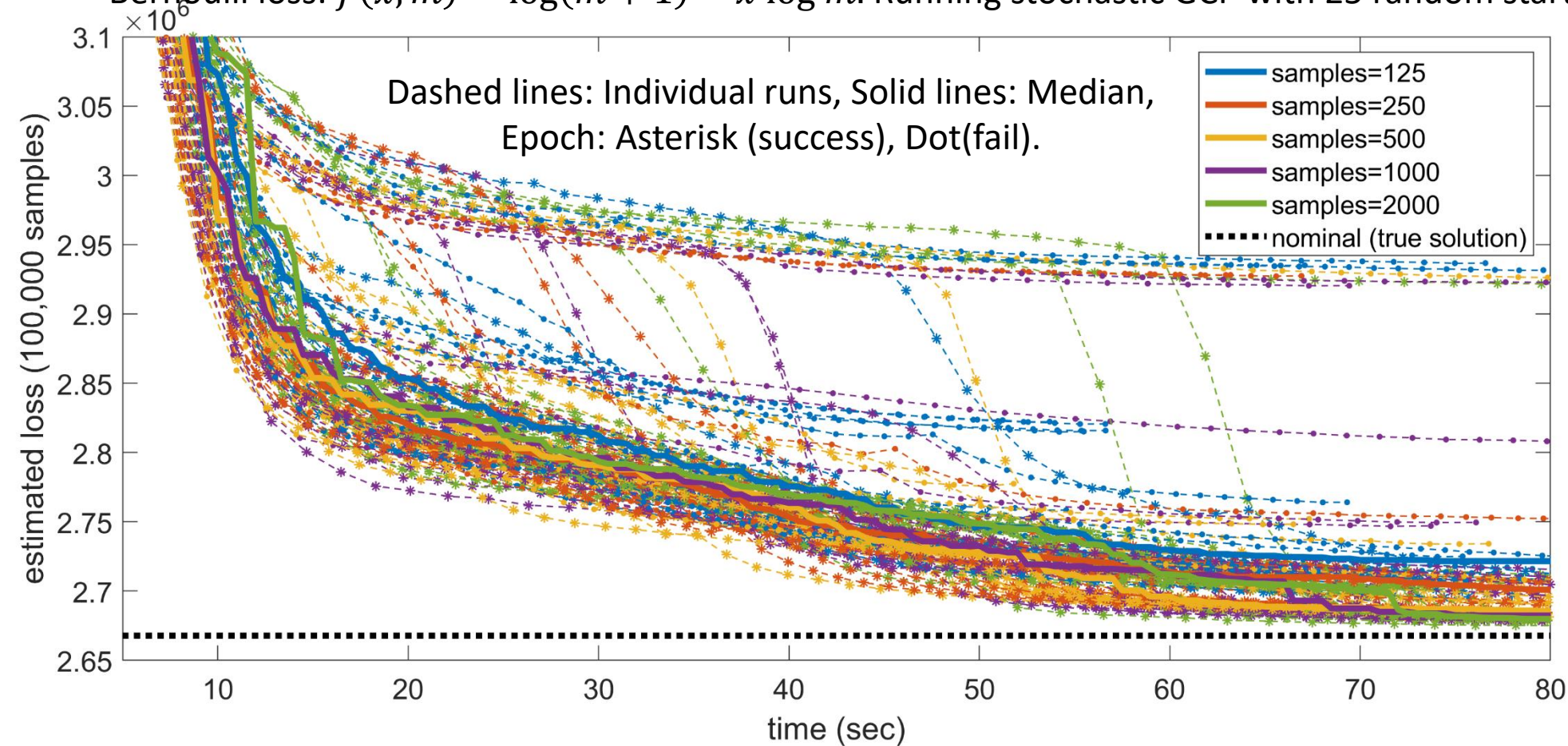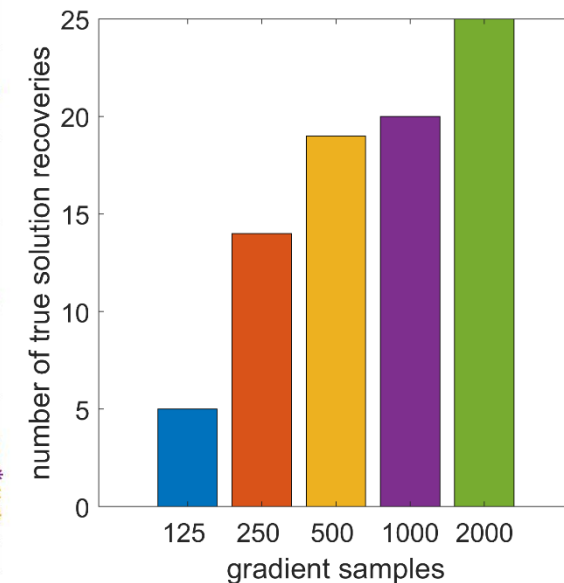Gamma loss: $f(x, m) = \frac{x}{m} + \log m$. Running stochastic GCP with 25 random starts and varying numbers of samples.



Dashed lines: Individual runs, Solid lines: Median, Epoch: Asterisk (success), Dot(fail).

Success at Recovering Underlying Generative Factors

# Stochastic vs. Non-Stochastic

$200 \times 150 \times 100 \times 50$ Tensor with low-rank ($r = 5$) structure based on Gamma distribution ($k = 1, \theta$ from model).

Gamma loss: $f(x, m) = \dfrac{x}{m} + \log m$. Running stochastic GCP with 25 random starts.



Each asterisk is an iteration.

Same as prior slide, but rescaled x-axis

Legend:
- samples = 125
- samples = 250
- samples = 500
- samples = 1000
- samples = 2000
- Non-stochastic
- nominal (true solution)

y-axis: estimated loss (100,000 samples), ×10⁷
x-axis: time (sec)

$200 \times 150 \times 100 \times 50$ Tensor with low-rank ($r = 5$) structure based on Bernoulli distribution (odds from model).
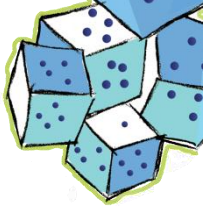Sparse tensor, less than 0.35% dense (~500K nonzeros).
Bernoulli loss: $f(x, m) = \log(m + 1) - x \log m$. Running stochastic GCP with 25 random starts, varying # of samples.
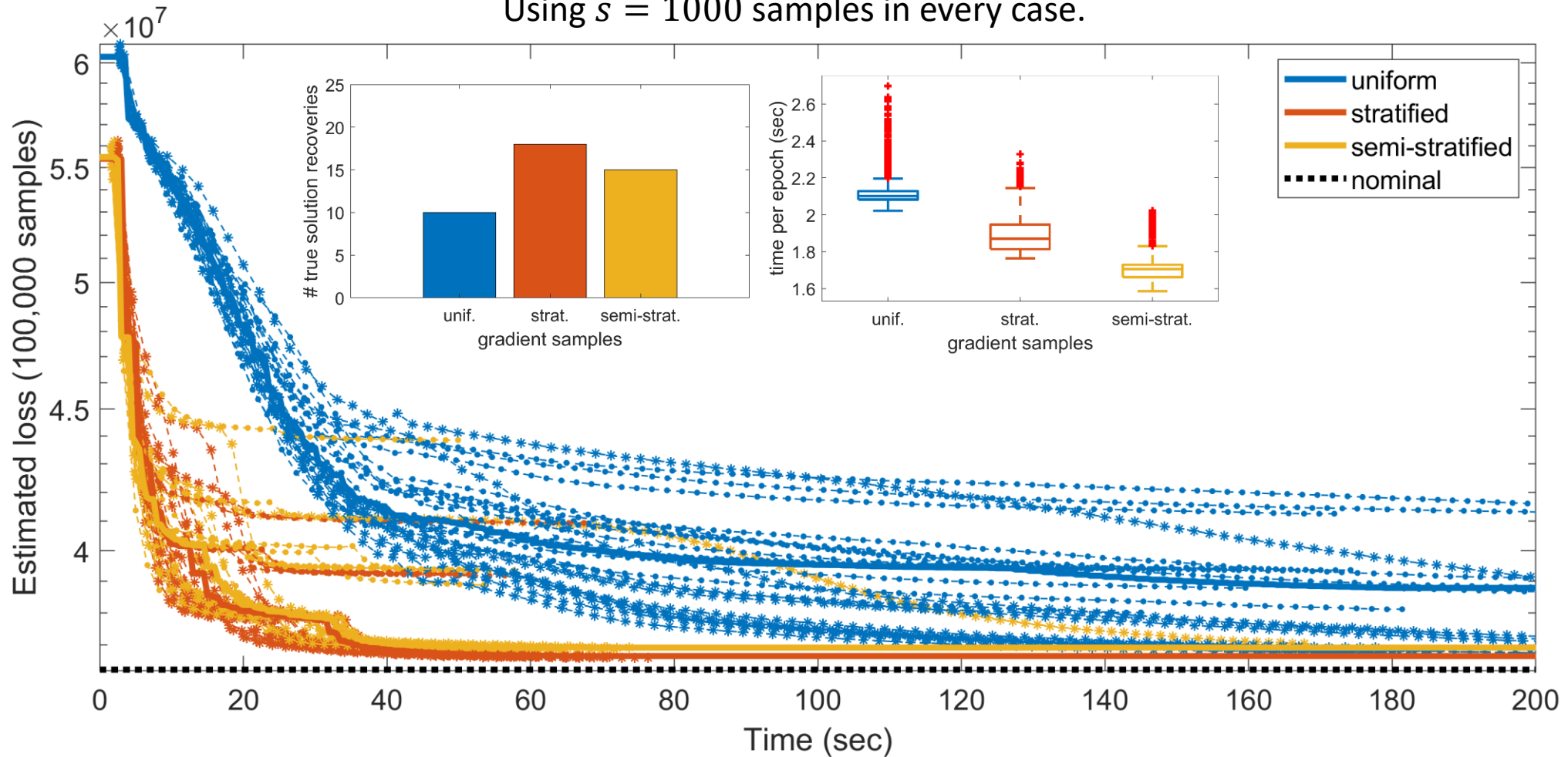


Dashed lines: Individual runs, Solid lines: Median,
Epoch: Asterisk (success), Dot(fail).

Legend:
- samples=125
- samples=250
- samples=500
- samples=1000
- samples=2000
- nominal (true solution)

Success at Recovering Underlying Generative Factors

# Uniform Sampling is Worse than Stratified for Sparse Tensors

Same set-up as binary experiments, but bigger tensor: $400 \times 300 \times 200 \times 100$, 0.38% dense (9M nonzeroes).
Using $s = 1000$ samples in every case.
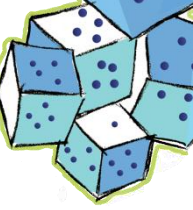
# Chicago Crime Data

- 4-way count tensor
  - 6,186 Days
  - 24 Hours of the Day
  - 77 Community Areas
  - 32 Crime Types
- Non-zeros: 5,330,673
  - Storage: 0.21GB for sparse tensor
- Distribution of entries
  - 0: 98.54%
  - 1: 1.33%
  - $\geq 2$: 0.12%
- Obtained from FROSTT (http://frostt.io/tensors/chicago-crime/)
- Data originally from Chicago Data Portal (https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-present/ijzp-q8t2)

GCP-Count
Rank = 10
$s = 6,319$

$$f(x, m) = m - x \log m$$



City of Chicago Community Areas and 'Sides'

# Application to Sparse Crime Binary Tensor (Semi-stratified results)
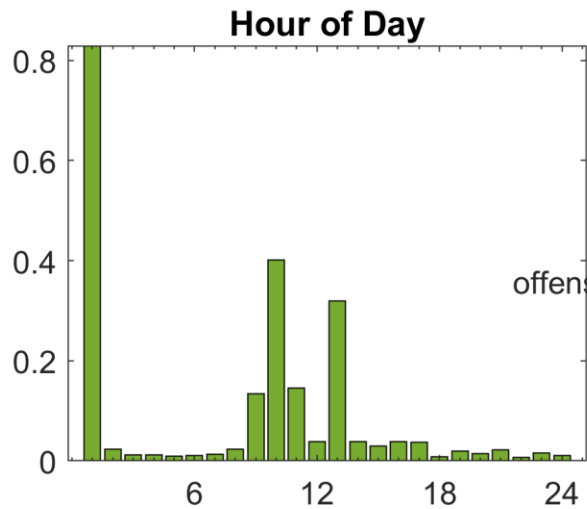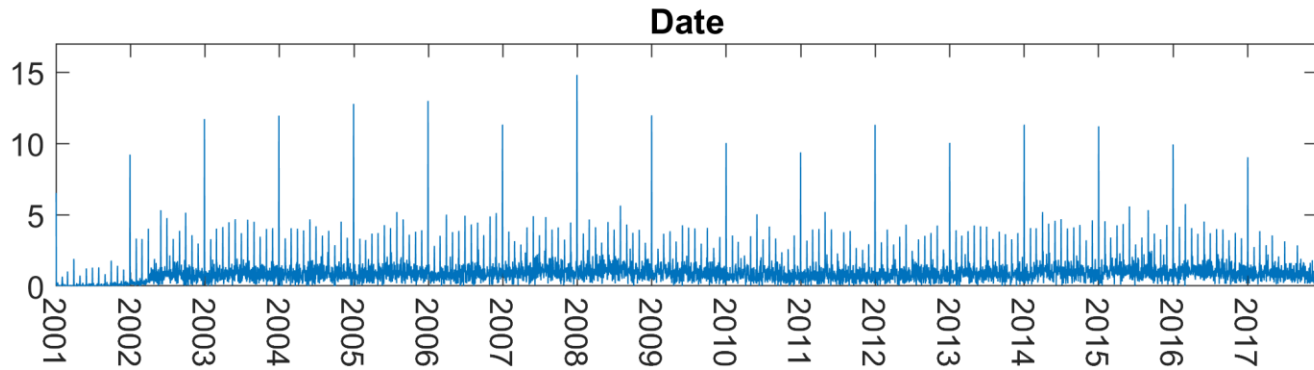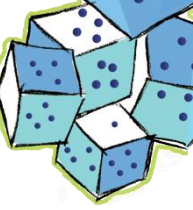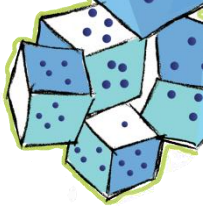
# Component #1

# Component #3

# Component #6

# Aside: Estimating Higher-Order Moments via Symmetric Tensor Factorization

Joint work with Sam Sherman, Notre Dame

Given a set of $p$ observations: $\mathbf{a}_i \in \mathbb{R}^n, i = 1, 2, \ldots, p$

First-order moment (mean): $\dfrac{1}{p} \sum_{i=1}^{p} \mathbf{a}_i$

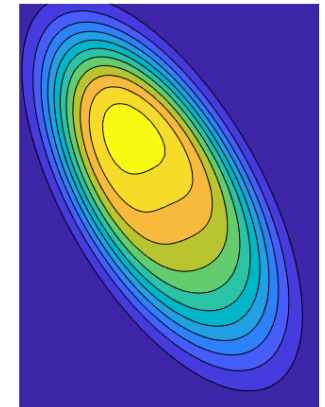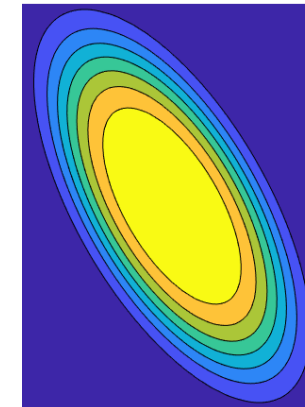Second-order moment: $\dfrac{1}{p} \sum_{i=1}^{p} \mathbf{a}_i \circ \mathbf{a}_i$

Third-order moment: $\dfrac{1}{p} \sum_{i=1}^{p} \mathbf{a}_i \circ \mathbf{a}_i \circ \mathbf{a}_i$

Fourth-order moment: $\dfrac{1}{p} \sum_{i=1}^{p} \mathbf{a}_i \circ \mathbf{a}_i \circ \mathbf{a}_i \circ \mathbf{a}_i$

We can compute low-rank ($r \ll p$) symmetric tensor estimated to higher-order moments…

$\dfrac{1}{r} \sum_{i=1}^{r} \mathbf{b}_i \circ \mathbf{b}_i \circ \mathbf{b}_i$

$\dfrac{1}{r} \sum_{i=1}^{r} \mathbf{c}_i \circ \mathbf{c}_i \circ \mathbf{c}_i \circ \mathbf{c}_i$

What are good applications, if any?

# References & Collaborators

*My department is hiring statisticians! Talk to me to learn more.*

- **Generalized CP (GCP) Tensor Decomposition** - D. Hong, T. G. Kolda, J. A. Duersch. **Generalized Canonical Polyadic Tensor Decomposition**, SIAM Review (to appear). http://arxiv.org/abs/1808.07452

- **Stochastic GCP** - D. Hong. T. G. Kolda. **Stochastic Gradients for Large-Scale Tensor Decomposition**, to appear on arXiv very soon!

- **Original mouse experiments** - A. H. Williams, T. H. Kim, F. Wang, S. Vyas, S. I. Ryu, K. V. Shenoy, M. Schnitzer, T. G. Kolda, S. Ganguli. **Unsupervised Discovery of Demixed, Low-dimensional Neural Dynamics across Multiple Timescales through Tensor Components Analysis**. *Neuron*, 98(6), 2018. https://doi.org/10.1016/j.neuron.2018.05.015

- **Poisson Tensor Factorization** - E. C. Chi, T. G. Kolda. **On Tensors, Sparsity, and Nonnegative Factorizations**. *SIAM Journal on Matrix Analysis and Applications*, 33(4), 2013. https://doi.org/10.1137/110859063

- **CP-APR Implementation** - S. Hansen, T. Plantenga, T. G. Kolda. **Newton-Based Optimization for Kullback-Leibler Nonnegative Tensor Factorizations**. *Optimization Methods and Software*, 30(5), 2015. https://doi.org/10.1080/10556788.2015.1009977

- **LDRD project team** - Cliff Anderson-Bergman (LLNL), Grey Ballard (Wake Forrest), Jed Duersch (SNL), Karen Devine (SNL), Srinivas Eswar (Georgia Tech), David Hong (Michigan), Jiajia Li (PNNL), Eric Phipps (SNL), Rich Vuduc (Georgia Tech), Jeff Young (Georgia Tech)

For more information and references: www.kolda.net