Large Data	Gen. Independence/Complexity		Image Summarization	Assay Selection	
11111	111111111	11111	11111	11	11

# Summarizing Large Data Sets

#### Jeffrey A. Bilmes

Professor Departments of Electrical Engineering & Computer Science and Engineering University of Washington, Seattle http://melodi.ee.washington.edu/~bilmes

#### Friday, March 27th, 2015

Large Data	Gen. Independence/Complexity	Doc Summarization	Data Summarization	Image Summarization	Assay Selection	
Out	line					



- Submodularity: Generalized Independence/Complexity
  - Document Summarization
  - Data Summarization
    - Speech Summarization
    - Selection in Statistical Machine Translation
    - Handwritten Digit Recognition
- Image Summarization
  - 6 Assay Selection

Large Data	Gen. Independence/Complexity	Doc Summarization	Data Summarization	Image Summarization	Assay Selection	End
Out	line					

### Large Data Sets

- 2 Submodularity: Generalized Independence/Complexity
- 3 Document Summarization
- Data Summarization
  - Speech Summarization
  - Selection in Statistical Machine Translation
  - Handwritten Digit Recognition
- 5 Image Summarization
- 6 Assay Selection



• Big Data is Really Big and Getting Even Bigger.



- Every day, we create 2.5 quintillion bytes (2.5 billion gigabytes) of data (source: IBM).
- 90% of the world's data has been created in the last two years.

# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End Big Data in Machine Learning Image Summarization Image Summarization

Statistics and Machine Learning

• "There's no data like more data", more samples reduces sampling error, higher statistical significance, and better *p* values.

# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End Big Data in Machine Learning Intervention Interve

- Statistics and Machine Learning
  - "There's no data like more data", more samples reduces sampling error, higher statistical significance, and better *p* values.
- Computational Consequences:
  - More expensive, computational resource demands (compute and storage), distributed implementations, more complicated

# Big Data in Machine Learning

Doc Summarization

Statistics and Machine Learning

Gen. Independence/Complexity

- "There's no data like more data", more samples reduces sampling error, higher statistical significance, and better *p* values.
- Computational Consequences:
  - More expensive, computational resource demands (compute and storage), distributed implementations, more complicated
  - Research opportunities to address new computational challenges





Image Summarization

Assav Selection

- systems programming, parallel and distributed computing, network topologies, efficient databases.
- Examples: map reduce, Hadoop, GraphLab, HaLoop, Greenplum, Asterix, Spark, SystemML, MLBase, Myria, etc.

Large Data

Large Data	Gen. Independence/Complexity		Image Summarization	Assay Selection	
Bigg	ger is Differe	ent			

Large Data	Gen. Independence/Complexity		Image Summarization	Assay Selection	
110111		11111			
Bigg	ger is Differe	ent			

• small (*n*-body)



Large Data	Gen. Independence/Complexity		Data Summarization	Image Summarization	Assay Selection	
		11111				
Bigg	ger is Differe	ent				

- small (*n*-body)
- medium (fluid dynamics, viscosity, compresibility),



Large Data	Gen. Independence/Complexity		Data Summarization	Image Summarization	Assay Selection	
		11111				
Bigg	ger is Differe	ent				

- small (*n*-body)
- medium (fluid dynamics, viscosity, compresibility),
- large (global weather systems, meteorology).



Same underlying molecular collision events!

Large Data	Gen. Independence/Complexity			Image Summarization	Assay Selection	
	111111111	11111	1111111111111			
Bigg	ger is Differe	ent				

Neurons

Bigger is Different	Large Data	Gen. Independence/Complexity		Image Summarization	Assay Selection	
Bigger is Different						
	Bigg	ger is Differe	ent			

Neurons

• small (neural spike trains, population coding)



Large Data	Gen. Independence/Complexity		Data Summarization	Image Summarization	Assay Selection	
Bigg	ger is Differe	ent				

Neurons

- small (neural spike trains, population coding)
- medium (intelligence, consciousness, psychology)





# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

Neurons

- small (neural spike trains, population coding)
- medium (intelligence, consciousness, psychology)
- large (society, social choice, wisdom of the crowd)







Same underlying electrical and chemical impulses.

Sum. Large Data — 3/27/2015

# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

 "More is Different", P.W. Anderson, 1972 (Nobel laureate). "The ability to reduce everything to simple fundamental laws does not imply the ability to start from those laws and reconstruct the universe."



# Bigger is Different

Gen. Independence/Complexity

Large Data

 "More is Different", P.W. Anderson, 1972 (Nobel laureate). "The ability to reduce everything to simple fundamental laws does not imply the ability to start from those laws and reconstruct the universe."

Doc Summarization

Summarization

Image Summarization

 "... alterations of being ... are not only the transition of one magnitude into another, but a transition from quantity into quality," Hegel, The Science of Logic, 1816



Assav Selection





• Hypothesis: extremely large data sets offer qualitatively different capabilities than small data sets.

# Big Data is Different Data: A Proposition

Doc Summarization

- Hypothesis: extremely large data sets offer qualitatively different capabilities than small data sets.
- Some Evidence: Image Completion (Hays & Efros, 2007)



Data Summarization

Image Summarization

Assav Selection

"our initial experiments ... on a dataset of ten thousand images were very discouraging. However, increasing the image collection to two million yielded a <u>qualitative</u> leap in performance"

Large Data

# Big Data is Different Data: A Proposition

Doc Summarization

- Hypothesis: extremely large data sets offer qualitatively different capabilities than small data sets.
- Some Evidence: Image Completion (Hays & Efros, 2007)



Data Summarization

Image Summarization

Assav Selection

"our initial experiments ... on a dataset of ten thousand images were very discouraging. However, increasing the image collection to two million yielded a <u>qualitative</u> leap in performance"

• Problem: Big data sets are big, unwieldy, computationally challenging, and highly redundant.

Large Data

# Big Data is Different Data: A Proposition

Doc Summarization

- Hypothesis: extremely large data sets offer qualitatively different capabilities than small data sets.
- Some Evidence: Image Completion (Hays & Efros, 2007)



Summarization

Image Summarization

Assav Selection

"our initial experiments ... on a dataset of ten thousand images were very discouraging. However, increasing the image collection to two million yielded a <u>qualitative</u> leap in performance"

- Problem: Big data sets are big, unwieldy, computationally challenging, and <u>highly redundant</u>.
- Research Quest: Can statistical predictions and actions be made <u>cost effectively</u> using the right small data?

Large Data

Large Data	Gen. Independence/Complexity	Data Summarization	Image Summarization	Assay Selection	
Out	line				

#### Large Data Sets

Submodularity: Generalized Independence/Complexity

3) Document Summarization

#### 4 Data Summarization

- Speech Summarization
- Selection in Statistical Machine Translation
- Handwritten Digit Recognition
- 5 Image Summarization
- 6 Assay Selection



### Sets and set functions

We are given a finite "ground" set of objects:



Also given a set function  $f : 2^V \to \mathbb{R}$  that valuates subsets  $A \subseteq V$ . Ex: f(V) = 6



### Sets and set functions

#### Subset $A \subseteq V$ of objects:

Also given a set function  $f : 2^V \to \mathbb{R}$  that valuates subsets  $A \subseteq V$ . Ex: f(A) = 1



#### Subset $B \subseteq V$ of objects:



Also given a set function  $f : 2^V \to \mathbb{R}$  that valuates subsets  $A \subseteq V$ . Ex: f(B) = 6

# Two Equivalent Submodular Definitions

#### Definition (submodular)

Large Data

Gen. Independence/Complexity

A function  $f : 2^V \to \mathbb{R}$  is submodular if for any  $A, B \subseteq V$ , we have that:

$$f(A) + f(B) \ge f(A \cup B) + f(A \cap B) \tag{1}$$

Data Summarization

Definition (submodular (diminishing returns))

A function  $f : 2^V \to \mathbb{R}$  is submodular if for any  $A \subseteq B \subset V$ , and  $v \in V \setminus B$ , we have that:

$$f(A \cup \{v\}) - f(A) \geq f(B \cup \{v\}) - f(B)$$

Incremental "value", "gain", or "cost" of v decreases (diminishes) as the context in which v is considered grows from A to B.

(2)

# Example: Number of Colors of Balls in Urns

Doc Summarization

• Consider an urn containing colored balls. Given a set S of balls, f(S) counts the number of distinct colors.

Summarization



Gen. Independence/Complexity

Large Data

Initial value: 2 (colors in urn). New value with added blue ball: 3



Initial value: 3 (colors in urn). New value with added blue ball: 3

Image Summarization

Assav Selection

• Submodularity: Incremental Value of Object Diminishes in a Larger Context (diminishing returns). Thus, *f* is submodular.

# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

• Unconstrained minimization and maximization:

$$\min_{X \subseteq V} f(X) \qquad (3) \qquad \max_{X \subseteq V} f(X)$$

 Knowing nothing about f, need 2<sup>n</sup> queries for any quality <sup>(a,b,c)</sup> assurance on candidate solution. Otherwise, solution can be unboundedly poor!!



(4)

# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

• Unconstrained minimization and maximization:

$$\min_{X \subseteq V} f(X) \qquad (3) \qquad \max_{X \subseteq V} f(X) \qquad (4)$$

 Knowing nothing about f, need 2<sup>n</sup> queries for any quality <sup>(a,b,c)</sup> assurance on candidate solution. Otherwise, solution can be unboundedly poor!!



• When f is submodular, however, Eq. (3) is polytime, and Eq. (4) is constant-factor approximable.



• Constrained case: interested only in a subset of subsets  $S \subseteq 2^V$ .

### Constrained Discrete Optimization

Doc Summarization

• Constrained case: interested only in a subset of subsets  $S \subseteq 2^V$ .

Data Summarization

• Ex: Bounded size  $S = \{S \subseteq V : |S| \le k\}$ , or given cost vector w and budget, bounded cost  $\{S \subseteq V : \sum_{s \in S} w(s) \le b\}$ .

Gen. Independence/Complexity

Large Data



Image Summarization

Assav Selection

# **Constrained Discrete Optimization**

Doc Summarization

• Constrained case: interested only in a subset of subsets  $S \subseteq 2^V$ .

Data Summarization

- Ex: Bounded size  $S = \{S \subseteq V : |S| \le k\}$ , or given cost vector w and budget, bounded cost  $\{S \subseteq V : \sum_{s \in S} w(s) \le b\}$ .
- Ex: feasible sets S as combinatorial objects.

Gen. Independence/Complexity

Large Data



Image Summarization

# Constrained Discrete Optimization

Doc Summarization

• Constrained case: interested only in a subset of subsets  $S \subseteq 2^V$ .

Data Summarization

- Ex: Bounded size  $S = \{S \subseteq V : |S| \le k\}$ , or given cost vector w and budget, bounded cost  $\{S \subseteq V : \sum_{s \in S} w(s) \le b\}$ .
- Ex: feasible sets & as combinatorial objects.

Gen. Independence/Complexity

• Ex: feasible sets S as sub-level sets of g,  $S = \{S \subseteq V : g(S) \le \alpha\}$ , sup-level sets  $S = \{S \subseteq V : g(S) \ge \alpha\}$ 



Image Summarization

Large Data

# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

#### • Constrained discrete optimization problems:

maximize	f(S)		minimize	f(S)	
subject to	<i>S</i> ∈ <i>S</i>	(5)	subject to	<i>S</i> ∈ S	(6)
$P \subset OV$	. L. C.	11 I C I			

where  $S \subseteq 2^{V}$  is the feasible set of sets.

### Doc Summarization Constrained Discrete Optimization

#### • Constrained discrete optimization problems:

maximize	f(S)		minimize	f(S)	
subject to	<i>S</i> ∈ S	(5)	subject to	<i>S</i> ∈ S	(6)

Data Summarization

Image Summarization

Assav Selection

Enc

where  $S \subseteq 2^V$  is the feasible set of sets.

• Fortunately, when f (and g) are submodular, these problems can often be solved with guarantees, often very efficiently! (Feige, Mirrokni & Vondrák 20XX; Goel, Karande, Tripathi & Wang; Svitkina & Fleischer 2010; Jegelka & Bilmes 2011, Iyer, Jegelka, & Bilmes 2013, Iyer & Bilmes 2014, and many others).

Large Data

# The Entropy Function

Gen. Independence/Complexity

• Given a joint distribution  $p(x_V)$  over |V| random variables:

$$f(A) \triangleq H(X_A) = -\sum_{x_A} p(x_A) \log p(x_A)$$
(7)

Image Summarization

Assav Selection

Enc

Summarization

with  $p(x_V)$  joint probability distribution over  $X_V$ .

Doc Summarization



Shannon's incredible 1948 paper stated that entropy is subadditive H(X<sub>A</sub>) + H(X<sub>B</sub>) ≤ H(A ∪ B). McGill 1954 stated that <u>further</u> conditioning reduces entropy, H(X<sub>A</sub>|X<sub>B</sub>) ≥ H(X<sub>A</sub>|X<sub>B</sub>, X<sub>C</sub>).
 This condition is identical to submodularity!

Large Data
## A Salmagundi of Submodularity

Doc Summarization

Gen. Independence/Complexity

Large Data

 Many diverse functions are submodular. E.g., set cover, vertex cover, edge cover, graph cut, bipartite neighborhoods, facility location, sums of weighted concave composed with additive functions, matrix rank, matroid rank, entropy, KL-divergence functions, quantum entropy, log determinant, spectral functions applied to Hermitian matrix, etc.

Data Summarization

Image Summarization

Assav Selection

### A Salmagundi of Submodularity

Gen. Independence/Complexity

Large Data

 Many diverse functions are submodular. E.g., set cover, vertex cover, edge cover, graph cut, bipartite neighborhoods, facility location, sums of weighted concave composed with additive functions, matrix rank, matroid rank, entropy, KL-divergence functions, quantum entropy, log determinant, spectral functions applied to Hermitian matrix, etc.

Summarization

Image Summarization

Assav Selection

Enc

- All submodular functions express a form of "abstract combinatorial independence" or "generalized complexity"
- Given submodular f,  $\exists$  a notion of "independence", i.e.,  $A \perp\!\!\!\perp B$ :

$$f(A \cup B) = f(A) + f(B), \tag{8}$$

• and a notion of "conditional independence", i.e.,  $A \perp\!\!\!\perp B | C$ :

$$f(A \cup B \cup C) + f(C) = f(A \cup C) + f(B \cup C)$$
(9)

### Submodularity as a Model for Real-World Data

Data Summarization

Doc Summarization

- Submodularity: useful model to valuate set of data objects.
- Meaning of valuations typically depend on if function is to be maximized or minimized.

Maximization diversity, coverage, span, irredundancy, and information.

Gen. Independence/Complexity

Large Data



Minimization cooperative costs, complexity, roughness, and irregularity.

Image Summarization

Assav Selection



## Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Image Summarization Assay Selection End Big Data Summarization Intervention Submodular Approaches

- We are given a set indexed by V
- Approach: 1) find a good function f : 2<sup>V</sup> → ℝ<sub>+</sub> that represents information in V. 2) Then optimize f to obtain a subset.
- Heuristic: design f by hand, hoping that f is a good proxy for the information within V. Acknowledge that f is a surrogate objective, guarantees are only in terms of f.
- 2) Alternatively: attempt to learn f, or some aspects of a good f, in some fashion based on training data.
  - We report on both kinds of results for summarizing large data sets below.

Large Data	Gen. Independence/Complexity	Doc Summarization	Data Summarization	Image Summarization	Assay Selection	
		11111				
Out	line					

2 Submodularity: Generalized Independence/Complexity

#### Document Summarization

#### Data Summarization

- Speech Summarization
- Selection in Statistical Machine Translation
- Handwritten Digit Recognition

### 5 Image Summarization

6 Assay Selection

### **End**

3



Given a large set of documents, summarize it down to the essential set of sentences.





- The summary on the left is a subset of the summary on the right.
- Consider adding a new (blue) sentence to each of the two summaries.
- Marginal benefit of adding the new (blue) sentence to the smaller (left) summary is no less than the marginal benefit of adding blue sentence to the larger (right) summary.

## Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End In Submodularity for document summarization?

Exists many unintentional uses of submodularity in NLP community.

- E.g., maximum marginal relevance (MMR) (Carbonell & Goldstein, 1998) has a diminishing returns property.
- Modified MMR (McDonald, 2007)
- Concept-based approaches (Filatova & Hatzivassiloglou 2004; Takamura & Okumura, 2009; Riedhammer et al., 2010; Qazvinian et al., 2010).
- Automatic evaluation of candidate summarizes are submodular: ROUGE-N (Lin 2004) and Pyramid (Nenkova & Passonneau, 2004).
- Both ROUGE-N and Pyramid are parameterized by good quality summarizes produced by humans, used only for evaluation.



- $\mathbf{w}^{\top} \mathbf{f}_t(\mathbf{y})$  is a convex combination of submodular functions.
- Mixture weights can be learnt via structured max-margin.

$$\begin{array}{ll} \underset{\mathbf{w} \geq 0, \xi_t}{\text{minimize}} & \frac{1}{T} \sum_{t} \xi_t + \frac{\lambda}{2} \|\mathbf{w}\|^2 \\ \text{subject to} & \mathbf{w}^\top \mathbf{f}_t(\mathbf{y}^{(t)}) \geq \max_{\mathbf{y} \in \mathcal{Y}_t} \left( \mathbf{w}^\top \mathbf{f}_t(\mathbf{y}) + \ell_t(\mathbf{y}) \right) - \xi_t, \forall t \\ & \xi_t \geq 0, \forall t. \end{array}$$
(10)

• Exponential set of constraints reduced to an embedded optimization problem, "inference."

Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection Internation Typical Results (DUC-06) Rouge-2: higher is better



Large Data	Gen. Independence/Complexity	Data Summarization	Image Summarization	Assay Selection	
Out	line				

2 Submodularity: Generalized Independence/Complexity

#### 3) Document Summarization

### Data Summarization

- Speech Summarization
- Selection in Statistical Machine Translation
- Handwritten Digit Recognition

### 5 Image Summarization

6 Assay Selection

Large Data	Gen. Independence/Complexity		Data Summarization	Image Summarization	Assay Selection	
		11111				

### As the data set size grow . . .

There is no data like more data

Large Data	Gen. Independence/Complexity	Doc Summarization	Data Summarization	Image Summarization	Assay Selection	End
		11111				
As t	he data set	size grov	N			

#### There is no data like more data $\Rightarrow$ more data is like no more data.

## Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

There is no data like more data  $\Rightarrow$  more data is like no more data.



Large Data	Gen. Independence/Complexity	Doc Summarization	Data Summarization	Image Summarization	Assay Selection	End
11111		1111				

### As the data set sizes grow ...





- Research question: Can statistical predictions be cost effective using small data?
- Using a submodular model that is easy to evaluate, can we produce useful and easy to obtain training data subsets that still perform well enough?

Large Data	Gen. Independence/Complexity	Data Summarization	Image Summarization	Assay Selection	
Out	line				

2 Submodularity: Generalized Independence/Complexity

3) Document Summarization

### Data Summarization

### Speech Summarization

- Selection in Statistical Machine Translation
- Handwritten Digit Recognition

### 5 Image Summarization

6 Assay Selection

### Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

### Speech Subset Selection

- Corpus Summarization: Given a large set of speech utterances (training corpus) V = {v<sub>1</sub>, v<sub>2</sub>, ..., v<sub>n</sub>}, choose a small subset A ⊆ V that is representative of V.
- Goal: training on summary should yield highest accuracy possible.
- Focus on <u>drastic</u> reductions in training set (one to two orders magnitude) to reduce model design time.

## Submodular Switchboard: GMM and DNN

Independence/Complexity

Large Data

	1%	5%	10%	20%	all
Rand	$52.1\pm1.5$	38.2±0.2	35.1±0.3	34.4±0.2	
HE (words)	49.6	36.5	34.8	N/A	
HE (3-phones)	47.5	37.6	34.2	N/A	21.0
SM (3-phones)	47.5	35.7	33.3	32.6	51.0

Data Summarization

Table : Word error rates, random (Rand), histogram-entropy (HE), the submodular (SM) system. Histogram-entropy results saturate after 10%.

	1%	5%	10%	20%	all
Rand	43.7±0.5	34.3±0.9	31.5±0.5	29.8±0.2	
HE (3-phones)	42.8	33.9	31.3	N/A	26.0
SM (3-phones)	41.1	31.8	29.3	28.2	

Table : Word error rates for DNN system.

Large Data	Gen. Independence/Complexity	Data Summarization	Image Summarization	Assay Selection	
Out	line				

2 Submodularity: Generalized Independence/Complexity

- B) Document Summarization
  - Data Summarization
    - Speech Summarization
    - Selection in Statistical Machine Translation
    - Handwritten Digit Recognition
- 5 Image Summarization
- 6 Assay Selection

### Data subset selection for machine translation

Doc Summarization

• Statistical Machine Translation (SMT): automatically translate from one human language to another.

Data Summarization

Image Summarization

Assav Selection

- Common problems in SMT: 1) test data is from a target domain while training data is mixed-domain; 2) phrase translation table, when based on all training data, can be massive.
- Solution: choose and then train using only a (domain-specific) subset of training data.
- Many previous approaches (e.g., *n*-gram overlap (Ittycheriah & Roukos, 2007), coverage of unseen *n*-grams (Eck et al. 2005), feature decay approach (Biçici & Yuret, 2011-2013)) are inadvertently submodular.
- Some (e.g., Moore & Lewis, 2010) are only modular.
- We approach directly using submodular functions.

Large Data



• V is set of sentences, U is a set of features (*n*-grams in our work).

## Feature based submodular functions

- V is set of sentences, U is a set of features (*n*-grams in our work).
- Feature-based submodular functions:

$$f(X) = \sum_{u \in U} w_u \phi_u(m_u(X))$$
(12)

where  $w_u > 0$  is a feature weight,  $m_u(X) = \sum_{x \in X} m_u(x)$  is a non-negative modular function specific to feature u,  $m_u(x)$  is a relevance score, a non-negative scalar score indicating the relevance of feature u in object x, and  $\phi_u$  is a u-specific non-negative non-decreasing concave function.

## Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection Results on Two Standard Translation Corpora

Method		Data Subset Sizes				
	10%	20%	30%	40%	100%	
		N	IST			
Random	0.3991	0.4142	0.4205	0.4220	0.4257	
Cross-entropy	0.4235	0.4292	0.4290	0.4292		
Submodular-5	0.4285	0.4356	0.4333	0.4324	0.4257	
Submodular-6	0.4302*	0.4334	0.4371*	0.4349		
		EUR	OPARL			
Random	0.2590	0.2652	0.2677	0.2697		
Cross-entropy	0.2639	0.2687	0.2704	0.2723	0.2651	
Submodular-5	0.2653	0.2727	0.2697	0.2720		
Submodular-6	0.2697*	0.2700	0.2740*	0.2723		

BLEU test-set scores (higher is better) for random, cross-entropy (standard baseline), and various submodular methods. Bold = significant over best Xent system.

	Gen. Independence/Complexity	Data Summarization	Image Summarization	Assay Selection	
Out	line	 			

2 Submodularity: Generalized Independence/Complexity

3 Document Summarization

#### Data Summarization

- Speech Summarization
- Selection in Statistical Machine Translation
- Handwritten Digit Recognition

### 5 Image Summarization

6 Assay Selection

### Handwritten Image Recognition

Doc Summarization

 Task: Train a machine learning system to classify hand-written digits correctly.

Data Summarization

 A standard data set MNIST (Lecun,'98).

Gen. Independence/Complexity

 Deep Neural Networks are state-of-the art on this kind of data, but they are very slow to train and greatly benefit from GPU hardware.

11/11//11// 2222222222222 5**3**33**3**333333 14444444444 55553555555 56666666666 +17777777777 

Image Summarization

Assav Selection

End

Large Data

### Image Recognition, Submodular Selection

Data Summarization

Doc Summarization



FASS = Filtered Active Submodular Selection

Assav Selection

End

$$\label{eq:US} \begin{split} \mathsf{US} &= \mathsf{Uncertainty} \\ \mathsf{Sampling} \end{split}$$

 $f_{\rm fs}$ ,  $f_{\rm fac}$ , and  $f_{\rm NN}$ are three submodular functions.

Large Data

Gen. Independence/Complexity

Large Data	Gen. Independence/Complexity		Image Summarization	Assay Selection	
11111		11111			
Out	line				

2 Submodularity: Generalized Independence/Complexity

- 3 Document Summarization
- Data Summarization
  - Speech Summarization
  - Selection in Statistical Machine Translation
  - Handwritten Digit Recognition

### Image Summarization

6 Assay Selection



### Modern Image collections

Many images, also that have a higher level gestalt than just a few.



Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

### Image Summarization

**Task:** Summarize collection of images by representative subset of the images

Applications:

- Summarizing your holiday pictures.
- Summarizing image search results
- Efficient browsing of image collections
- Video frame summarization









### Image Summarization - Data Collection

### **Data Statistics**

- 14 image collections with 100 pictures each
- $\bullet \sim 400$  human summaries for every image collection, via Amazon Turk, about 5500 summaries total!

Example collections:





Large Data	Gen. Independence/Complexity		Data Summarization	Image Summarization	Assay Selection	
		1111		111		

### Image Summarization

### Super-Pixel Based V-Rouge

Whole collection:



3 best summaries:



### 3 medium summaries:



### 3 worst summaries:





#### Typical Results - Learnt mixture using Max-Margin



Max of Learned Mixture

Average Pruned Human

Large Data	Gen. Independence/Complexity	Data Summarization	Image Summarization	Assay Selection	
Out	line				

2 Submodularity: Generalized Independence/Complexity

- 3 Document Summarization
- 4 Data Summarization
  - Speech Summarization
  - Selection in Statistical Machine Translation
  - Handwritten Digit Recognition

### Image Summarization

6 Assay Selection

### ) End

# Large Data Gen. Independence/Complexity Doc Summarization Data Summarization Image Summarization Assay Selection End

- Functional genomics: improve understanding of gene function (e.g., transcription, translation, and forms of interactions).
- Genomic assays measure a DNA activity in cell-type specific nucleus. E.g., transcription, transcription factor binding, DNA accessibility, covalent modification of the histone proteins, etc.
- Cost prohibitive: The cost of genomics assays limits their application. Ideally, to fully characterize a cell type, one desires every possible assay type, but too costly.



### Assay Selection: Typical Results

Doc Summarization

- Fortunately, assay types are often redundant.
- Goal: choose a subset of assay types such that remaining assay types can be deduced.

Data Summarization

 Below are some preliminary results on summarizing genomic assays, quality of summaries measured in a variety of ways.

metric	whole data set	averaged random	facility location	saturated coverage
annotation prediction	0.434	$0.270 {\pm} 0.039$	0.319	0.330
assay prediction (SVR)	0.690	$0.977 {\pm} 0.028$	0.898	0.943
assay prediction (LR)	0.728	$1.032{\pm}0.025$	0.969	0.975
function prediction (F-measure)	0.0489	$0.0029 {\pm} 0.0025$	0.0065	0.0081
function prediction (AUC-PR)	0.076	$0.027{\pm}0.009$	0.038	0.042

Large Data

Assav Selection
Large Data	Gen. Independence/Complexity		Data Summarization	Image Summarization	Assay Selection	End
		11111				
Out	line					

## Large Data Sets

2 Submodularity: Generalized Independence/Complexity

3 Document Summarization

#### 4 Data Summarization

- Speech Summarization
- Selection in Statistical Machine Translation
- Handwritten Digit Recognition

## 5 Image Summarization

6 Assay Selection

# 7 End

Large Data	Gen. Independence/Complexity		Image Summarization	Assay Selection	End
		1111			
Ack	nowledgmen	ts			

Thanks to the following amazing people (former & current students, and current colleagues):

Mukund Narasimhan, Hui Lin, Andrew Guillory, Stefanie Jegelka, Sebastian Tschiatschek, Kai Wei, Yuzong Liu, Rishabh Iyer, Jennifer Gillenwater, Yoshinobu Kawahara, Katrin Kirchhoff, Carlos Guestrin, & Bill Noble.



# The End: Thank you!

# $f(A) + f(B) \ge f(A \cup B) + f(A \cap B)$ $= f(A_r) + 2f(C) + f(B_r) = f(A_r) + f(C) + f(B_r) = f(A \cap B)$ $+ 0 \ge 0 + 0 + 0 = 0$