

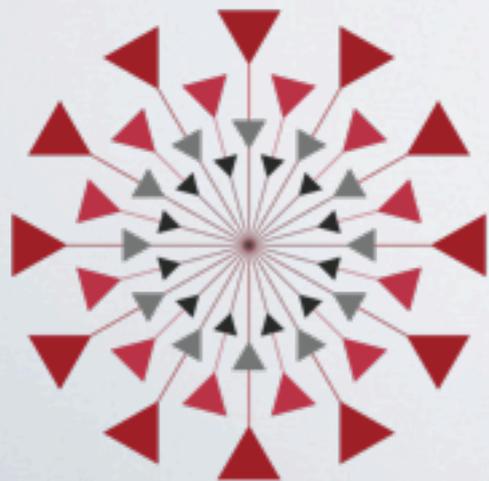
# BUILDING PROBABILISTIC MODELS AROUND DETERMINISTIC OPTIMIZATION PROCEDURES

Ryan Prescott Adams

School of Engineering and Applied Sciences  
Harvard University

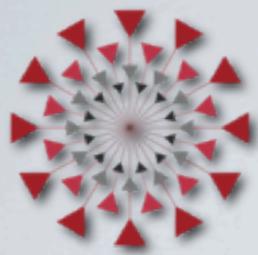
Joint work with Danny Tarlow and Rich Zemel

<http://hips.seas.harvard.edu>



**HARVARD**  
**INTELLIGENT**  
**PROBABILISTIC**  
**SYSTEMS**

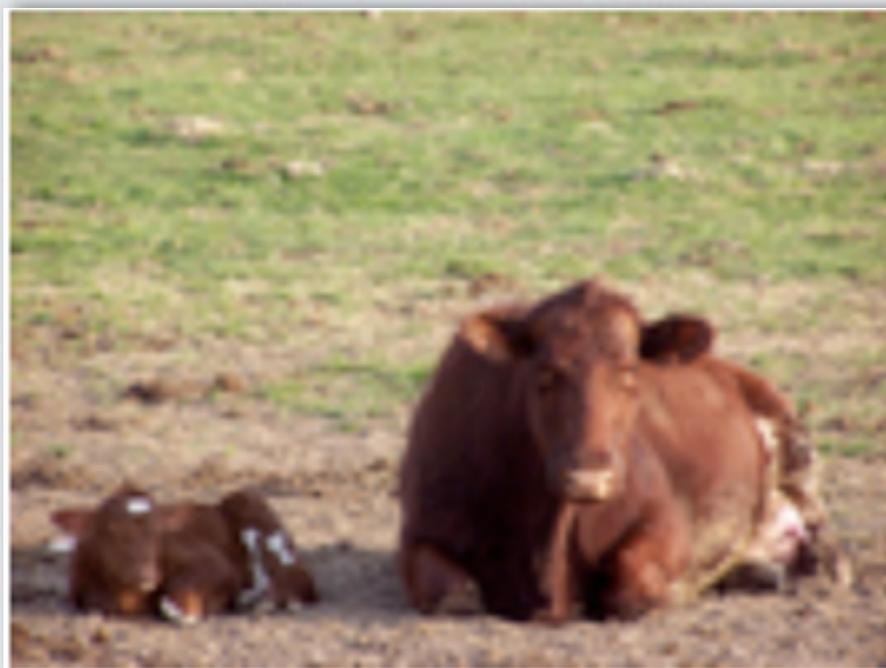




# STRUCTURED PREDICTION

---

- ▶ **Objective:** Build discriminative models that map from features to high-dimensional, structured discrete spaces.
- ▶ **Example: Image Segmentation**

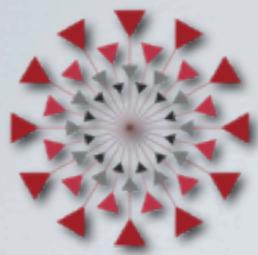


$$x \in \mathcal{X}$$

$$f(x; \theta) \rightarrow y$$



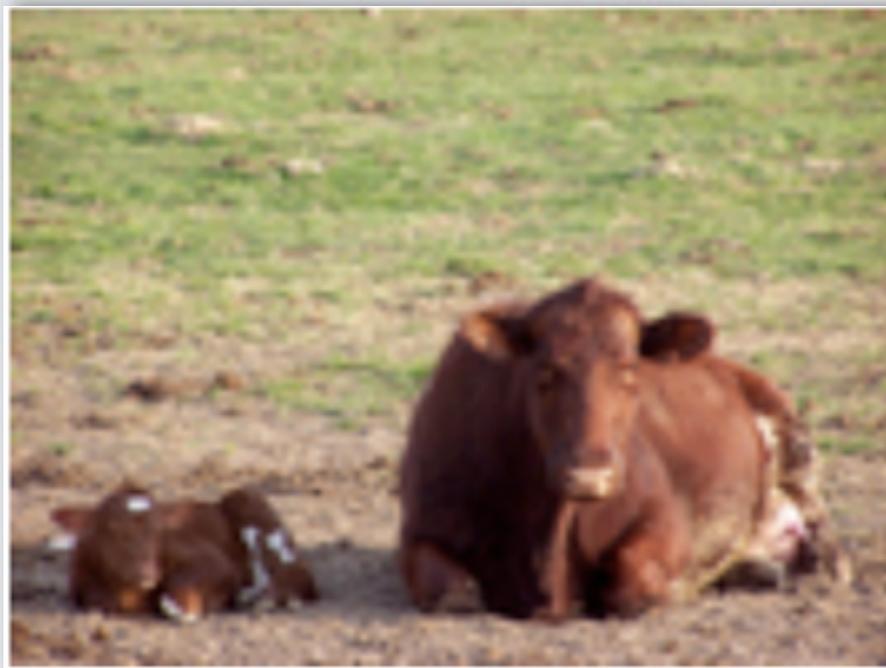
$$y \in \{1, 2, \dots\}^K$$



# STRUCTURED PREDICTION

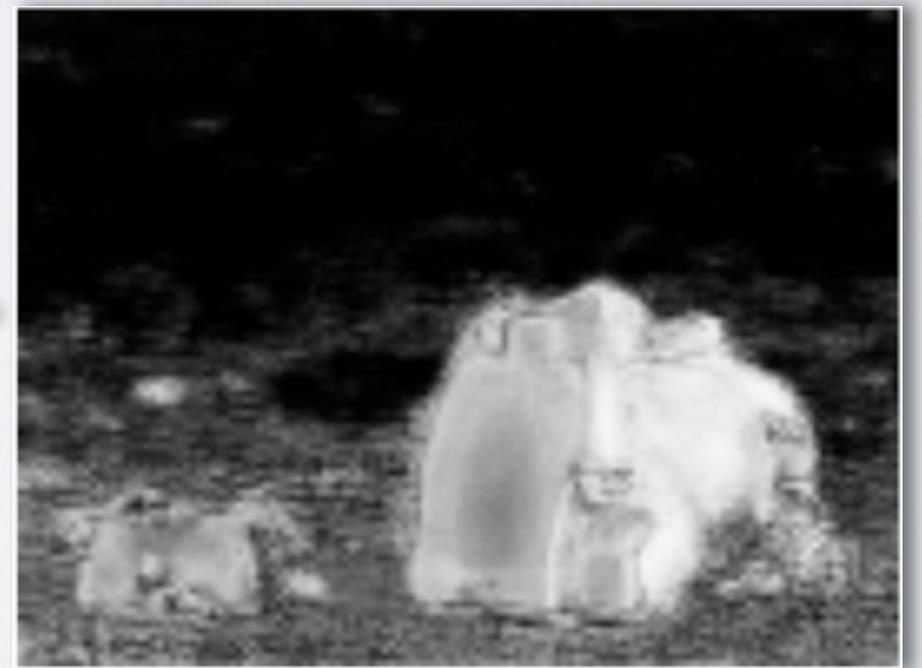
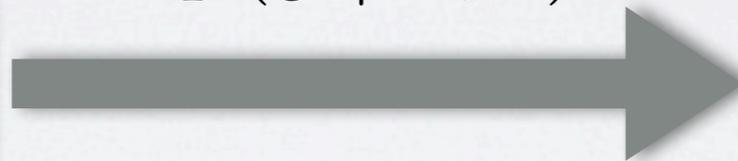
---

- ▶ **Objective:** Build discriminative models that map from features to high-dimensional, structured discrete spaces.
- ▶ **Example: Image Segmentation**

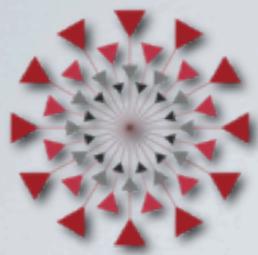


$$\mathbf{x} \in \mathcal{X}$$

$$p(\mathbf{y} \mid \mathbf{x}, \theta)$$



$$\mathbf{y} \in \{1, 2, \dots\}^K$$

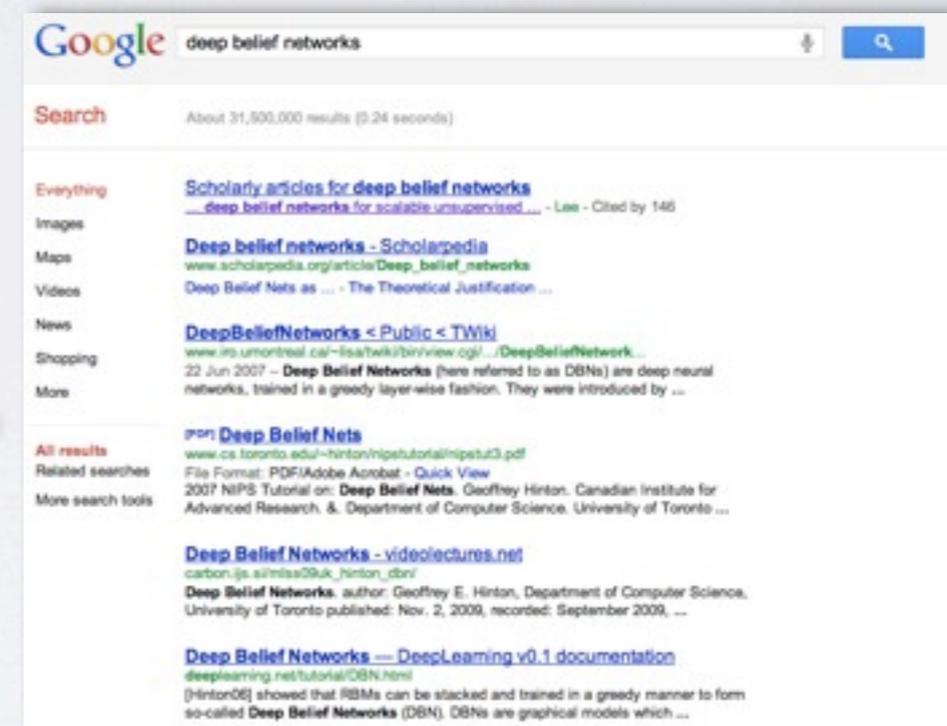
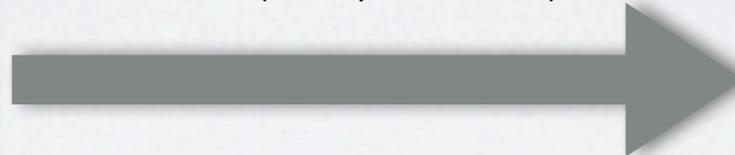


# STRUCTURED PREDICTION

- ▶ **Objective:** Build discriminative models that map from features to high-dimensional, structured discrete spaces.
- ▶ **Example: Ranking Documents**

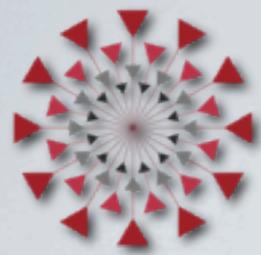


$$p(\mathbf{y} \mid \mathbf{x}, \theta)$$



$$\mathbf{x} \in \mathcal{X}$$

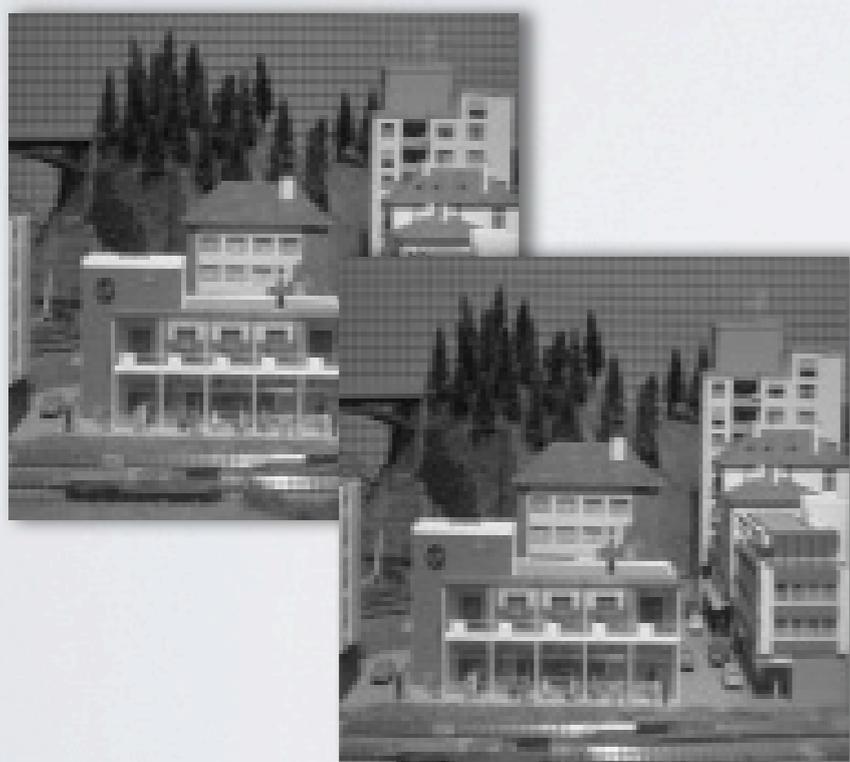
$$\mathbf{y} \in \{1, 2, \dots\}^K$$



# STRUCTURED PREDICTION

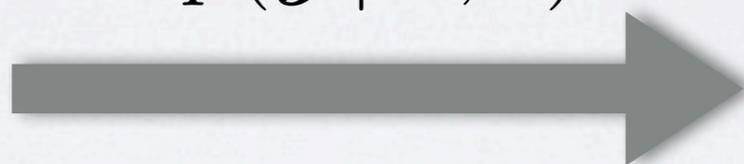
---

- ▶ **Objective:** Build discriminative models that map from features to high-dimensional, structured discrete spaces.
- ▶ **Example: Matching**

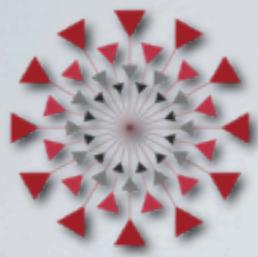


$$x \in \mathcal{X}$$

$$p(\mathbf{y} \mid x, \theta)$$



$$\mathbf{y} \in \{1, 2, \dots\}^K$$



# CONDITIONAL RANDOM FIELDS

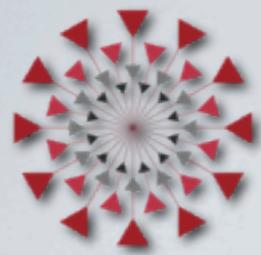
---

- ▶ A very common approach to structured prediction is to construct a feature-dependent energy.
- ▶ The resulting Gibbs distribution gives a probabilistic model for the discrete space:

$$p(\mathbf{y} \mid \mathbf{x}, \theta) = \frac{1}{\mathcal{Z}(\mathbf{x}, \theta)} \exp \{ -E(\mathbf{x}, \mathbf{y}; \theta) \}$$

$$E(\mathbf{x}, \mathbf{y}; \theta) = \theta^\top \psi(\mathbf{x}, \mathbf{y})$$

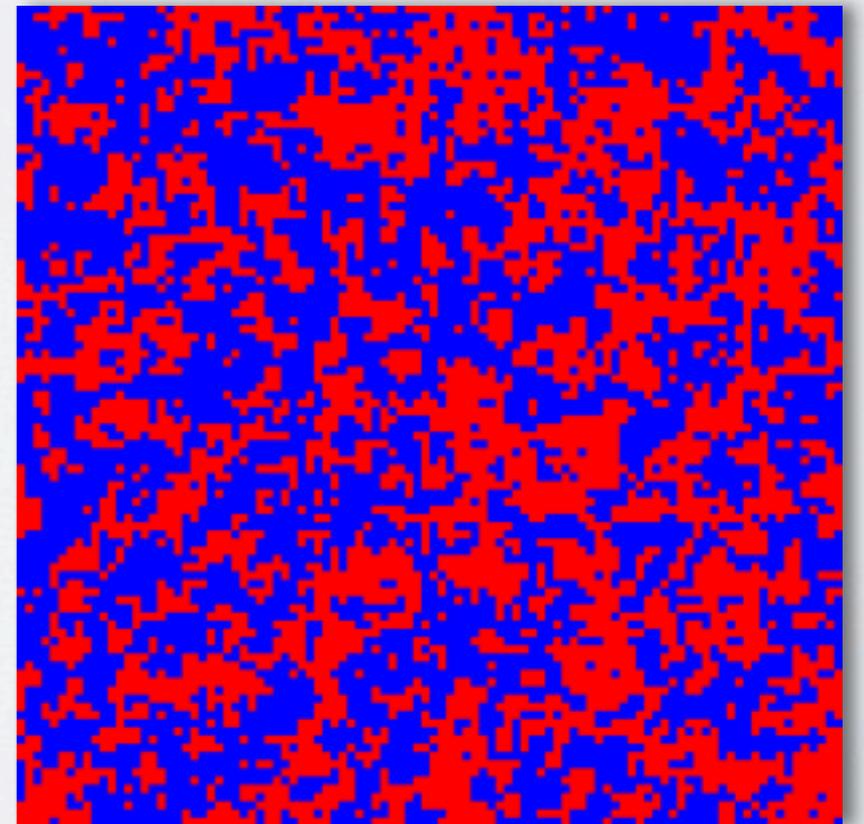
- ▶ A unique partition function for each training example!

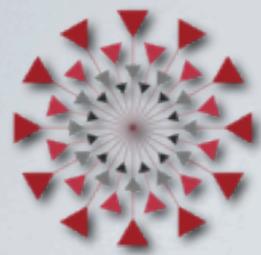


# A COMPUTATIONAL DISCREPANCY

---

- ▶ In many MRFs, finding the mode is easy even when finding the partition function is difficult. That is, one may be able to maximize in polynomial time over a set that may be exponentially large.
- ▶ **Example: Ising model with positive weights**
- ▶ Partition function is a sum over exponentially many configurations.
- ▶ Graph cuts finds the most probable configuration in polynomial time.



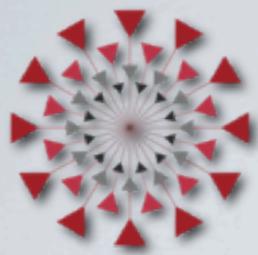


# A COMPUTATIONAL DISCREPANCY

---

- ▶ In many MRFs, finding the mode is easy even when finding the partition function is difficult. That is, one may be able to maximize in polynomial time over a set that may be exponentially large.
- ▶ **Example: Bipartite matching**
- ▶ Normalization is  $\#P$ -complete.
- ▶ The Hungarian algorithm finds the most probable match in cubic time.





# A COMPUTATIONAL DISCREPANCY

---

- ▶ In many MRFs, finding the mode is easy even when finding the partition function is difficult. That is, one may be able to maximize in polynomial time over a set that may be exponentially large.

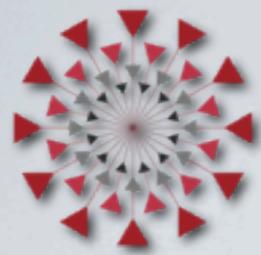
- ▶ **Example: Shortest paths**

- ▶ Normalization is #P-complete

- ▶ Dijkstra's algorithm finds the most probable solution in

$$O(|E| + |V| \log |V|)$$

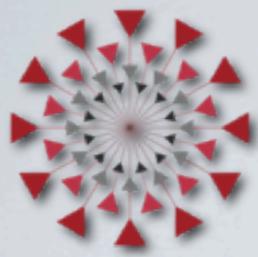




# A SIMPLE GENERATIVE MODEL

---

- ▶ Basic idea: Represent the predictive uncertainty in the choice of energy, rather than in the induced MRF.
- ▶ Use the efficient mode-finding algorithm as an explicit component of the data generation procedure.
- ▶ Use the features to induce a distribution on the parameters of an objective function, sample them and then optimize to generate the observed data.
- ▶ **Randomized Optimum Models (RandOMs)** are a generalization of Perturb-and-MAP [Papandreou & Yuille, 2011]



# RANDOMIZED OPTIMUM MODEL

---

- ▶ Define a family of objectives on the label space:

$$f_{\theta}(\mathbf{y}) : \mathcal{Y} \rightarrow \mathbb{R}$$

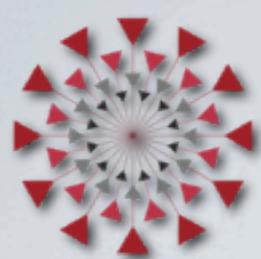
- ▶ Features of data:  $\phi(\mathbf{x})$

- ▶ Parameters of distribution over objectives:  $\gamma$

- ▶ Draw a random objective:  $\theta \sim p(\theta \mid \phi(\mathbf{x}), \gamma)$

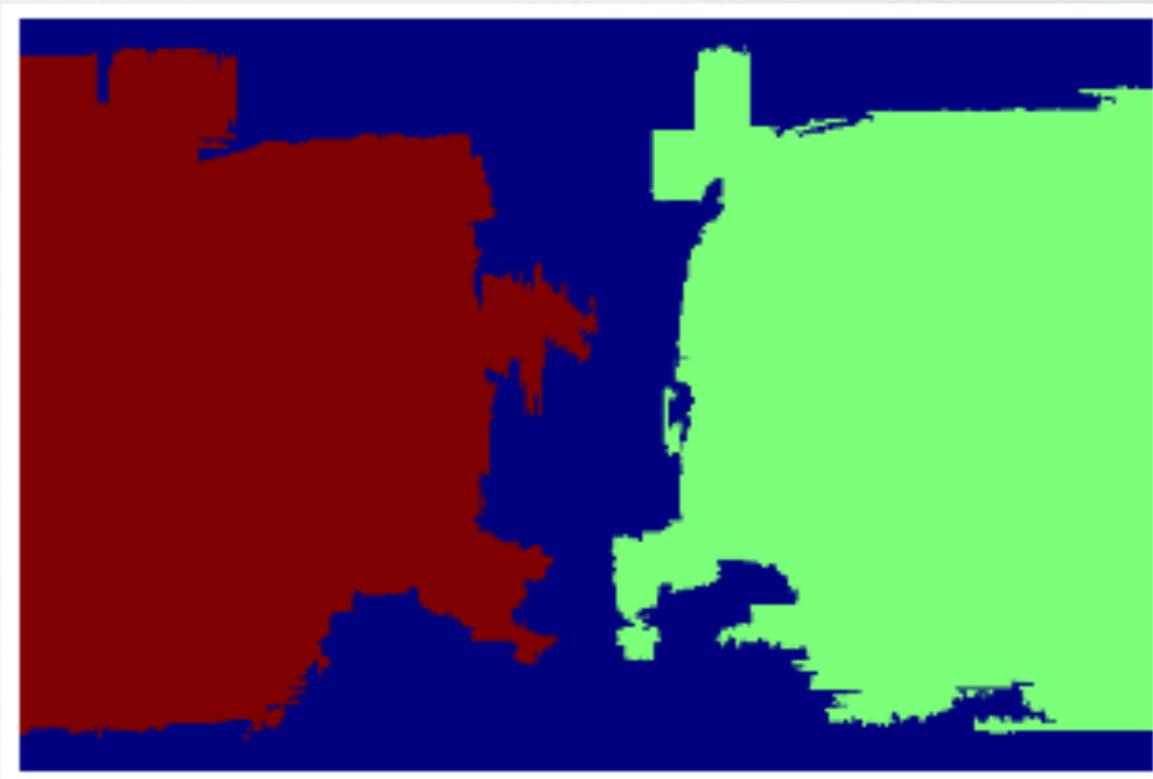
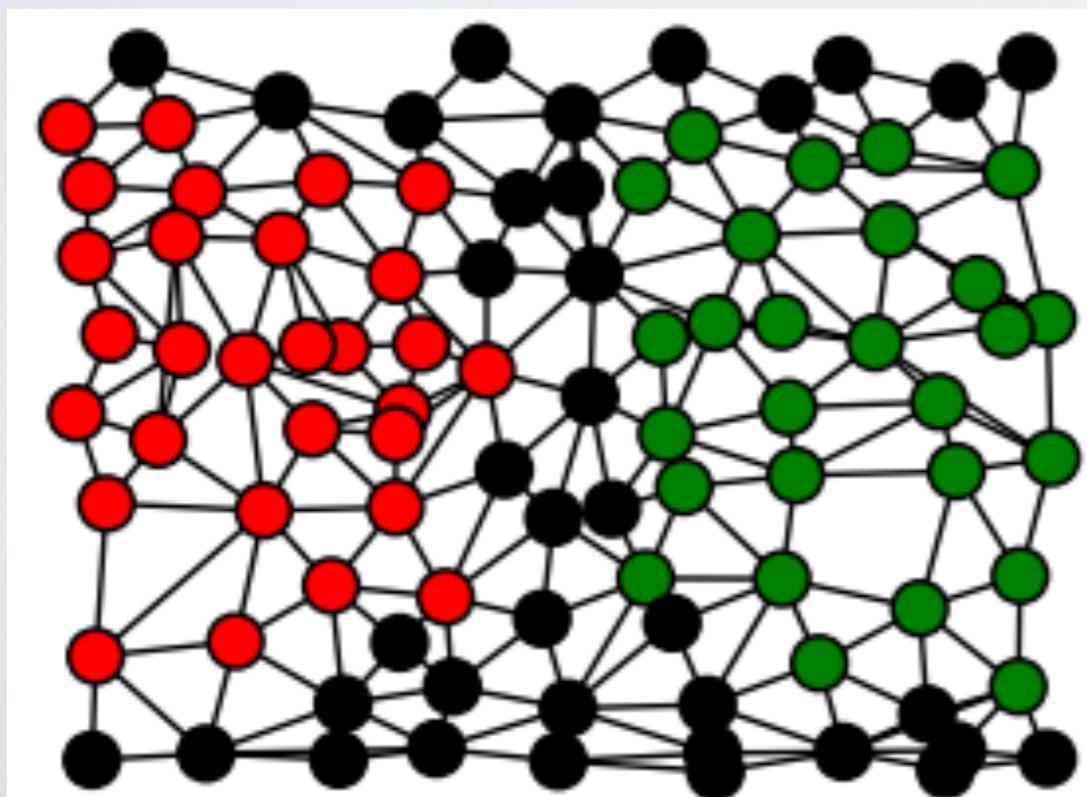
- ▶ Data are the minimizing configuration:

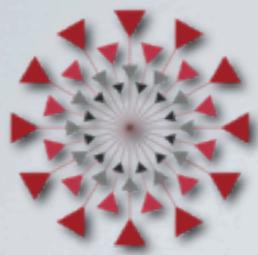
$$\mathbf{y} = \underset{\mathbf{y}'}{\operatorname{argmin}} f_{\theta}(\mathbf{y}')$$



# NON-MRF EXAMPLE: CONNECTED COMPONENTS

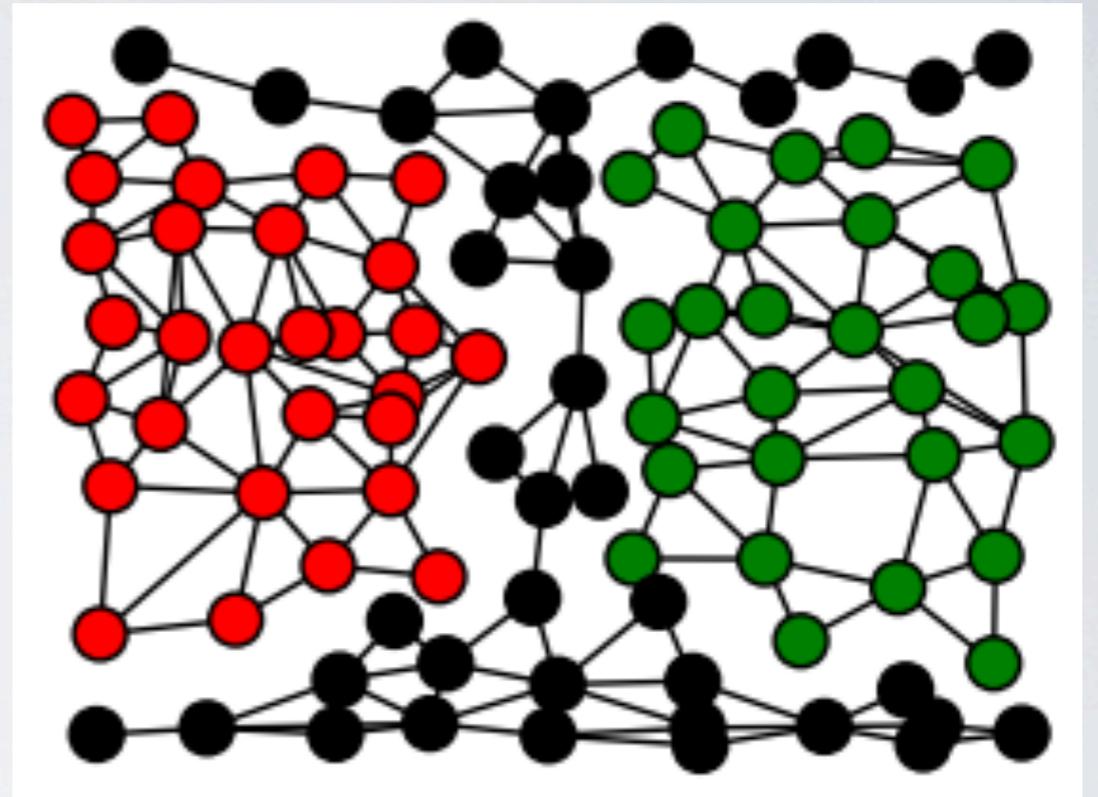
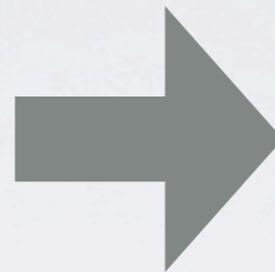
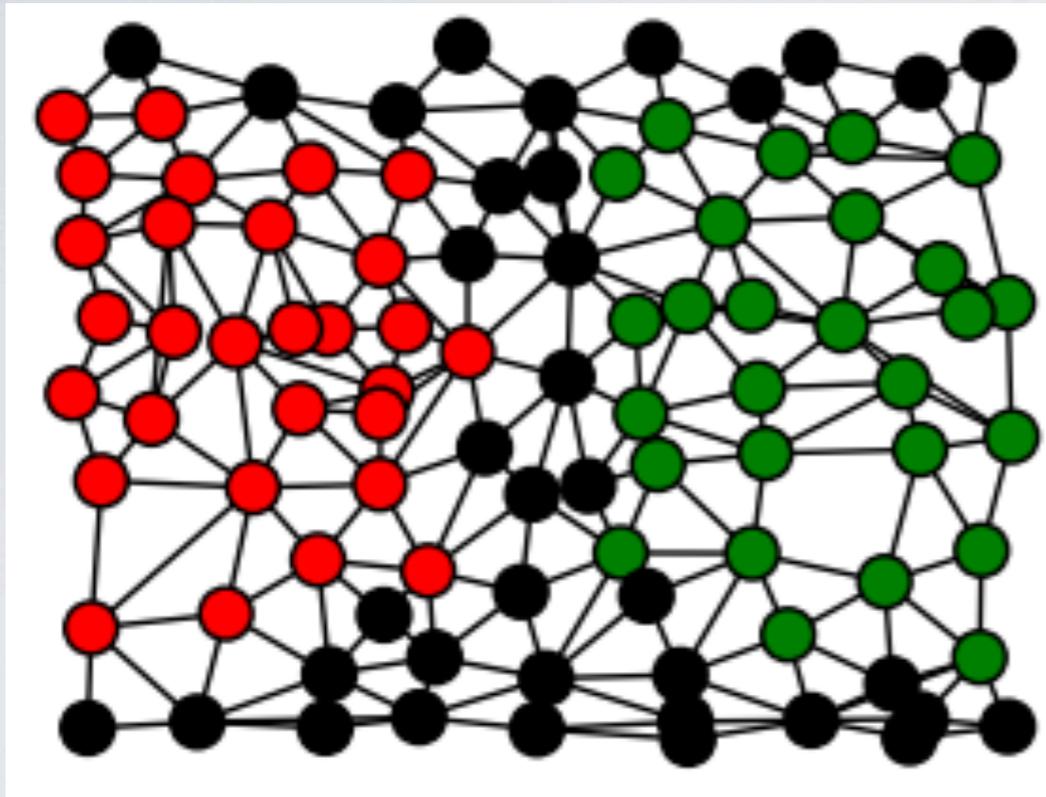
---



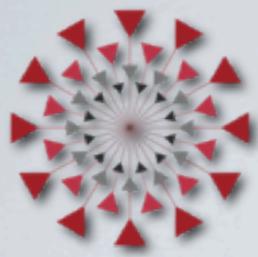


# NON-MRF EXAMPLE: CONNECTED COMPONENTS

---



- ▶ Draw random edge weights.
- ▶ Remove edges beneath threshold.
- ▶ Provides random multi-class segmentation.



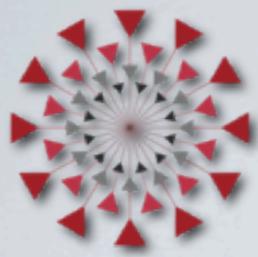
# LEARNING RANDOMIZED OPTIMUM MODELS

---

- ▶ We learn the maximum likelihood estimates of  $\gamma$ .
- ▶ Here, the MLE parameters are the ones that place the greatest mass on the objective functions that result in the observed data.
- ▶ Challenge: there are many objectives which are consistent with any given datum.

$$F(\theta) = \operatorname{argmin}_{\mathbf{y}'} f_{\theta}(\mathbf{y}')$$

$$\theta \in F^{-1}(\mathbf{y}) \Rightarrow F(\theta) = \mathbf{y}$$



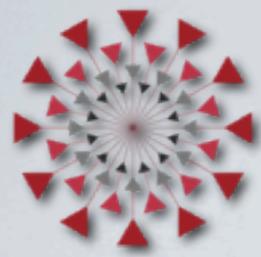
# LEARNING RANDOMIZED OPTIMUM MODELS

---

$$F(\theta) = \operatorname{argmin}_{\mathbf{y}'} f_{\theta}(\mathbf{y}')$$

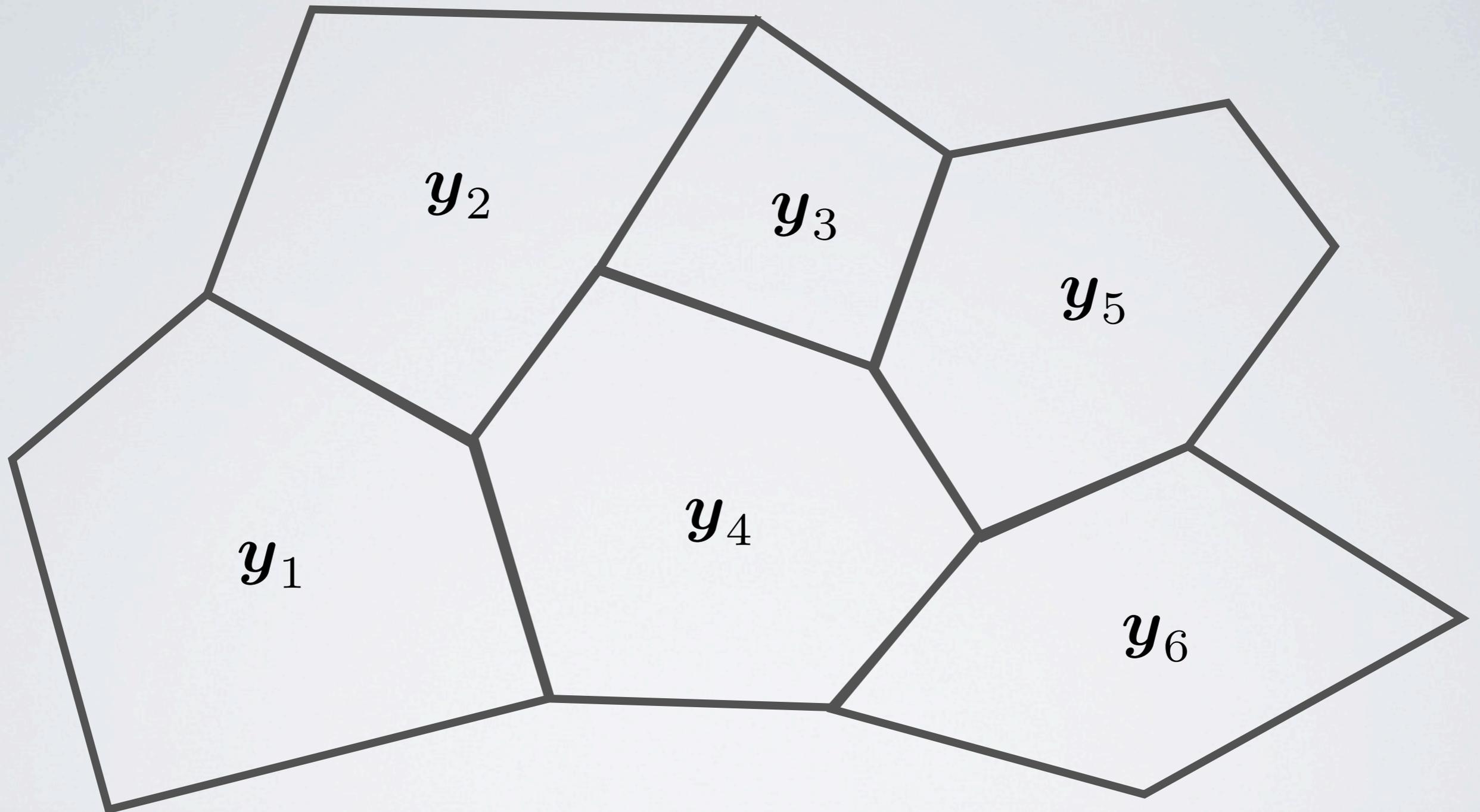
$$\theta \in F^{-1}(\mathbf{y}) \Rightarrow F(\theta) = \mathbf{y}$$

- ▶ The likelihood of  $\theta$  is uniform over  $F^{-1}(\mathbf{y})$ .
- ▶ We focus on cases where the sets  $F^{-1}(\mathbf{y})$  are simple, e.g., convex, star-convex, or easily projected upon.

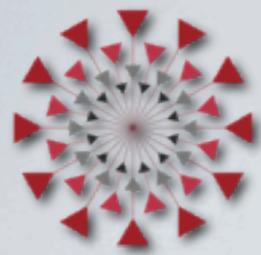


# THE GEOMETRY OF INVERSE SETS

---

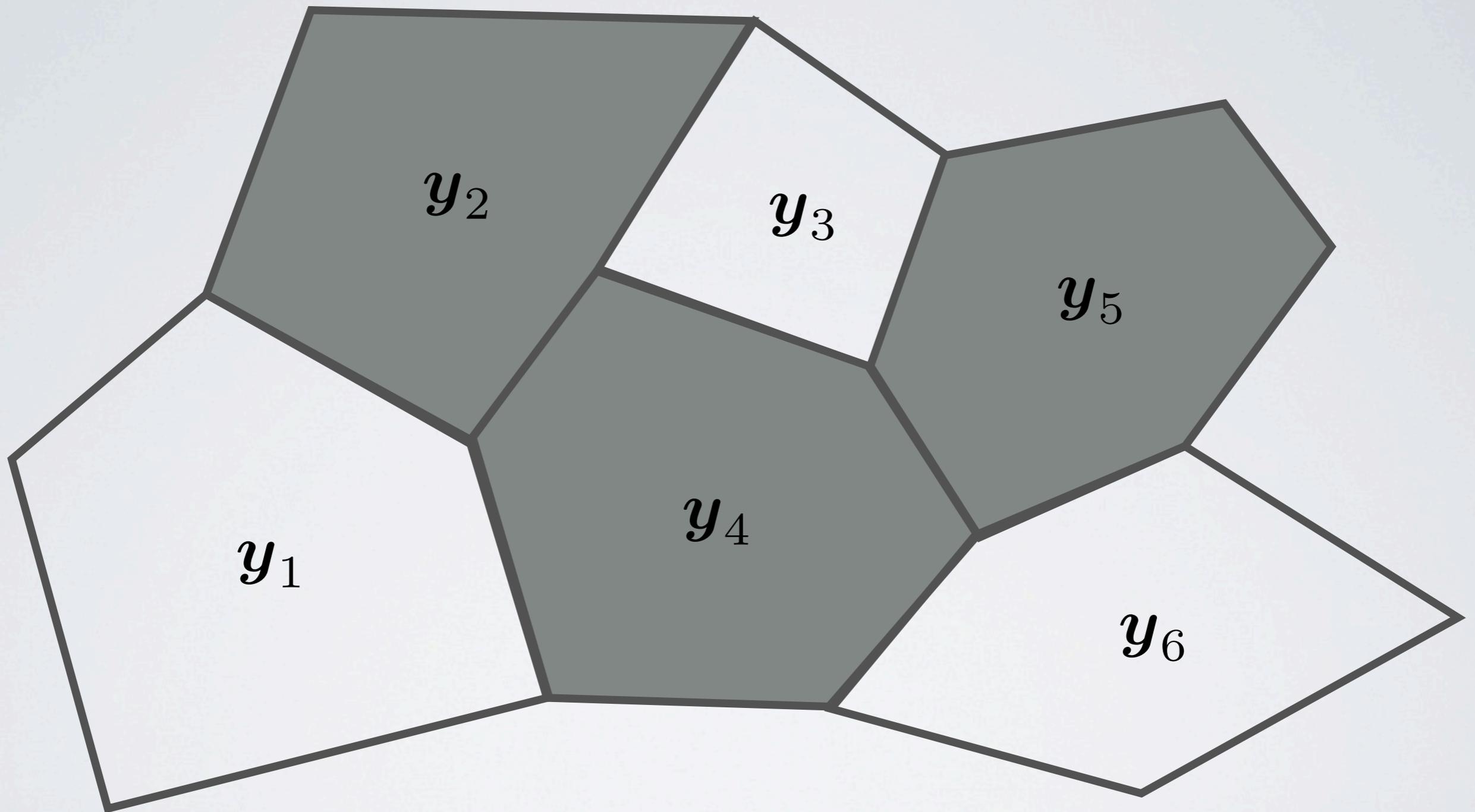


Equivalence classes of objective functions,  
leading to identical solutions.

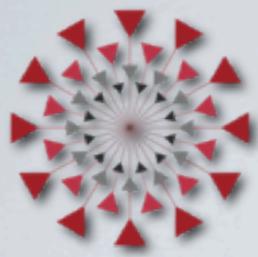


# THE GEOMETRY OF INVERSE SETS

---

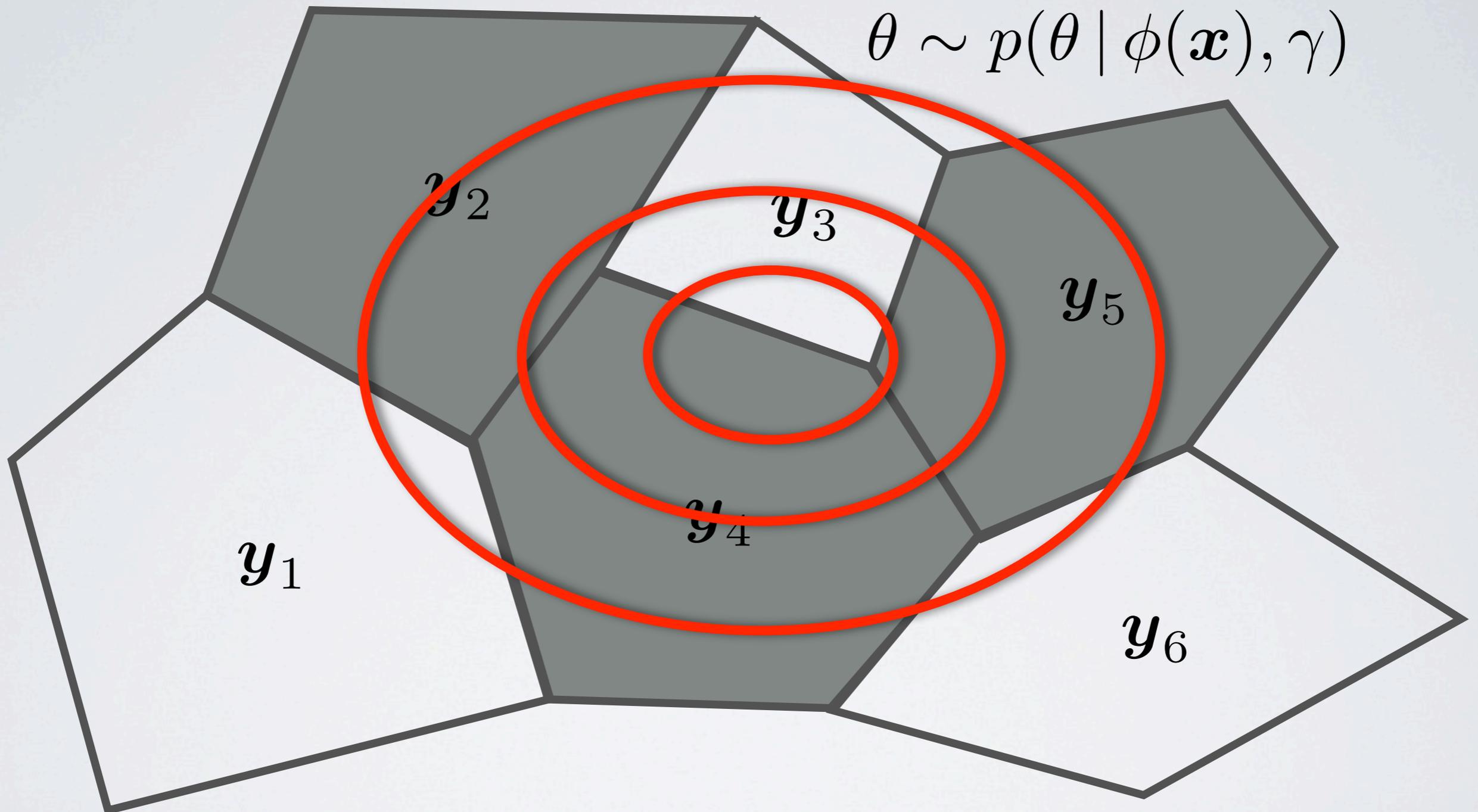


A subset of the equivalence classes are associated with observed data.

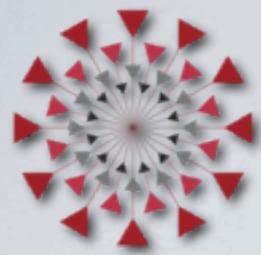


# THE GEOMETRY OF INVERSE SETS

---



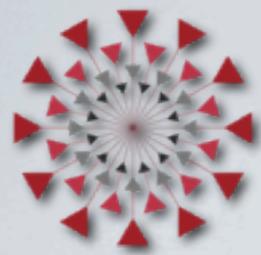
Find the distribution that places the most mass within the observed classes.



# MONTE CARLO EXPECTATION MAXIMIZATION

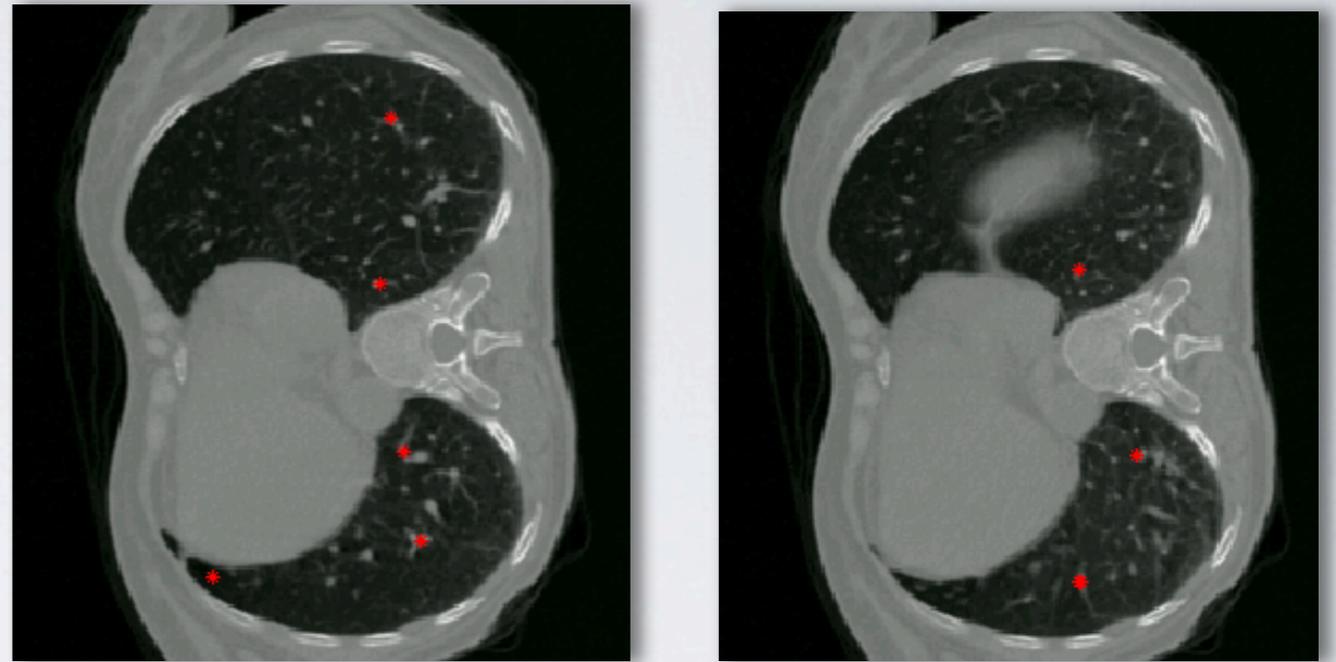
---

- ▶ We don't know which objective function is the "true" one for any given data example.
- ▶ We would like to integrate out the corresponding objective parameter  $\theta$ , to maximize the likelihood.
- ▶ Monte Carlo EM uses MCMC to integrate out the latent location in the data-consistent volume.
- ▶ The simple shape of the inverse set enables this to mix well.
- ▶ Dynamic variants of, e.g., graph cuts lead to efficient inside/outside oracles for slice sampling.

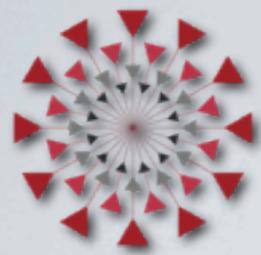


# LUNG CT RESULTS

- ▶ Volumetric CT of lungs in 10 patients.
- ▶ 30 landmarks each
- ▶ Find feature-driven matches.



Alg	Avg. Score	Avg. Prob	Avg. Log Prob
Max-margin	91.0%	.400	0 (6)
PM(1)	95.0%	.538	-1.52 (0)
PM(10)	94.9%	.543	-1.54 (0)
MCEM(1)	98.8%	.837	-1.67 (2)
MCEM(5)	96.8%	.568	-0.62 (1)
MCEM(10)	96.7%	.544	-0.56 (1)



# A NEW PROBABILISTIC MODELING TOOL

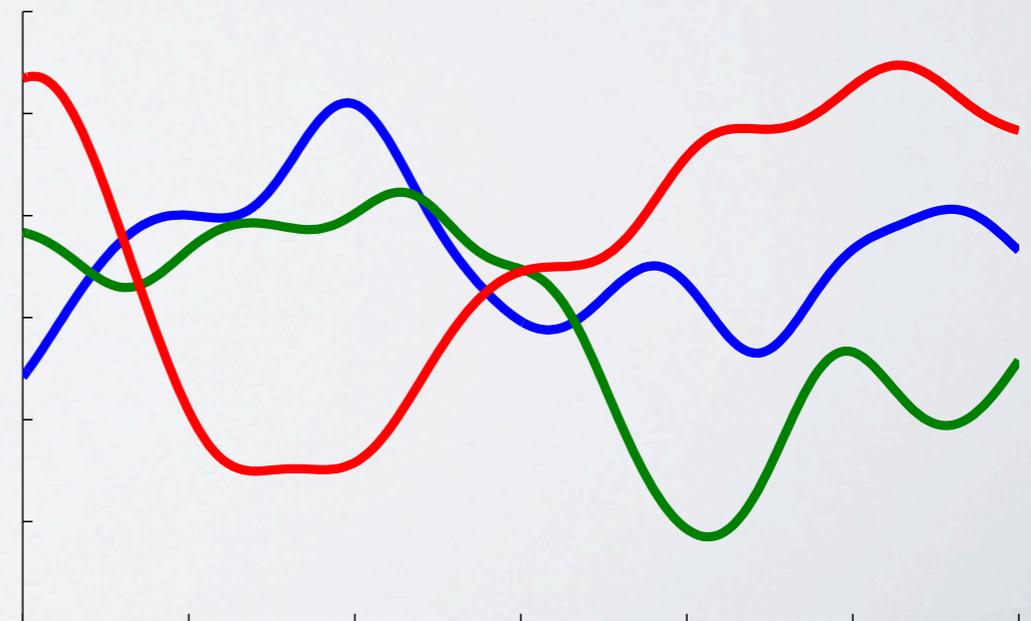
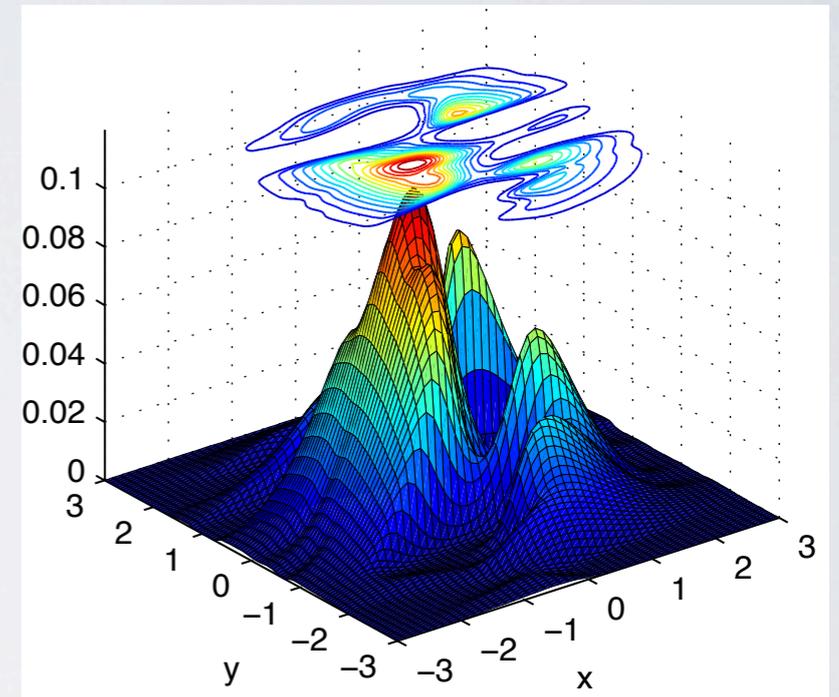
▶ We can import flexible tools for supervised and unsupervised modeling from continuous spaces and use them for structured prediction:

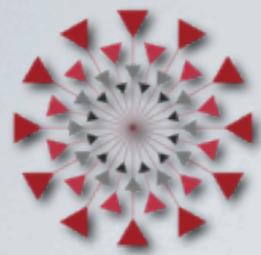
▶ Gaussian process regression

▶ Latent linear models

▶ Dirichlet process mixtures

▶ Latent feature models





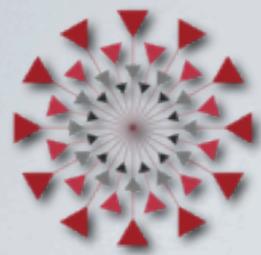
# UNCERTAINTY IN IMAGE SEGMENTATION

---

- ▶ When making predictions in low-level vision, we often want to represent our uncertainty.
- ▶ Acquiring marginals from MRFs can be expensive.
- ▶ Kohli & Torr (2008) suggested softmax min-marginals, with energy  $E(\mathbf{y})$  and  $y_d \in \{1, 2, \dots, K\}$ :

$$\Phi_d(k) = \min_{\mathbf{y}|y_d=k} E(\mathbf{y}) \quad p(y_d = k) \approx \frac{\exp\{-\Phi_d(k)\}}{\sum_{k'=1}^K \exp\{-\Phi_d(k')\}}$$

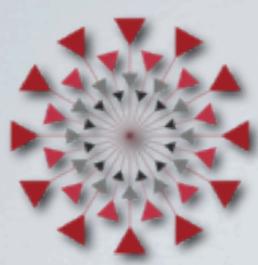
- ▶ Easy to compute, but may make very poor predictions.



# GENERATING DATA VIA MIN-MARGINALS

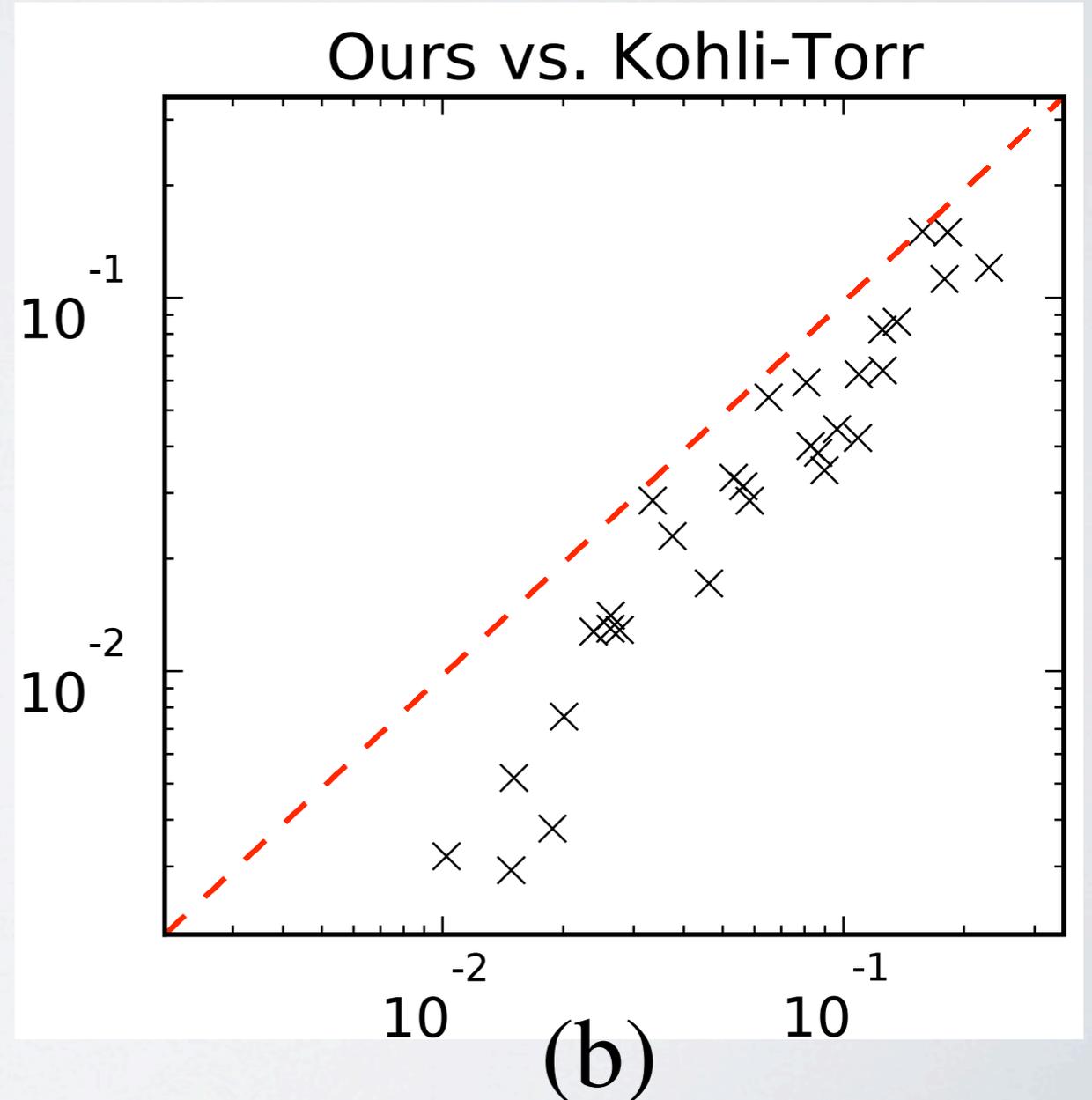
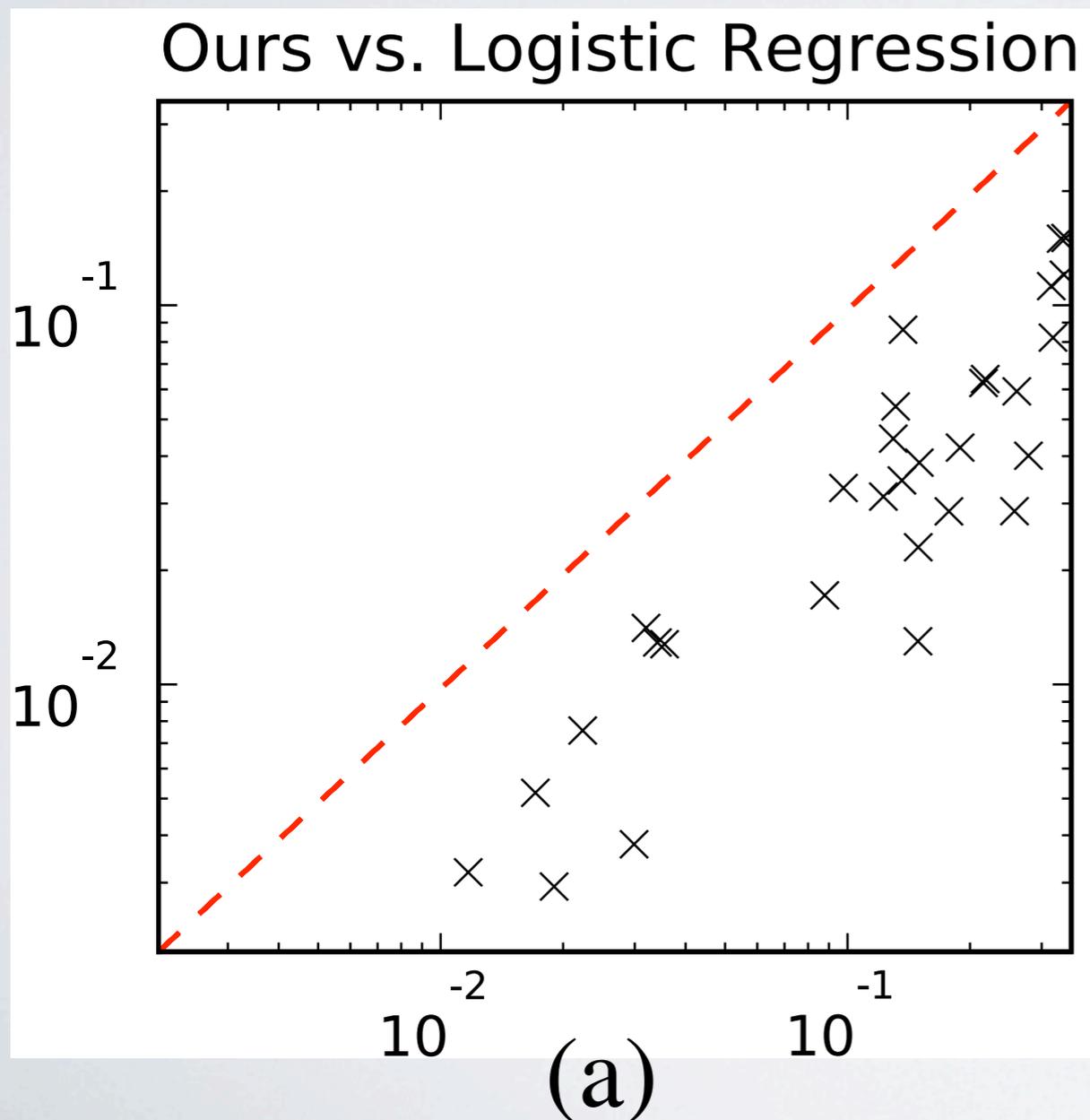
---

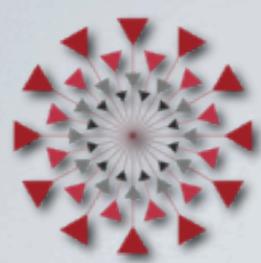
- ▶ As in RandOMs, why bother with the MRF at all?
- ▶ We take Kohli & Torr (2008) seriously as providing a probabilistic model for segmentation.
- ▶ **Effectively a coupled parameterization of a large collection of independent multinomials.**
- ▶ It is possible to compute the gradient of the log likelihood in terms of the underlying CRF parameters.
- ▶ The objective is to train for achieving well-calibrated test-time marginal predictions.



# EVALUATING TRAINING PROCEDURES

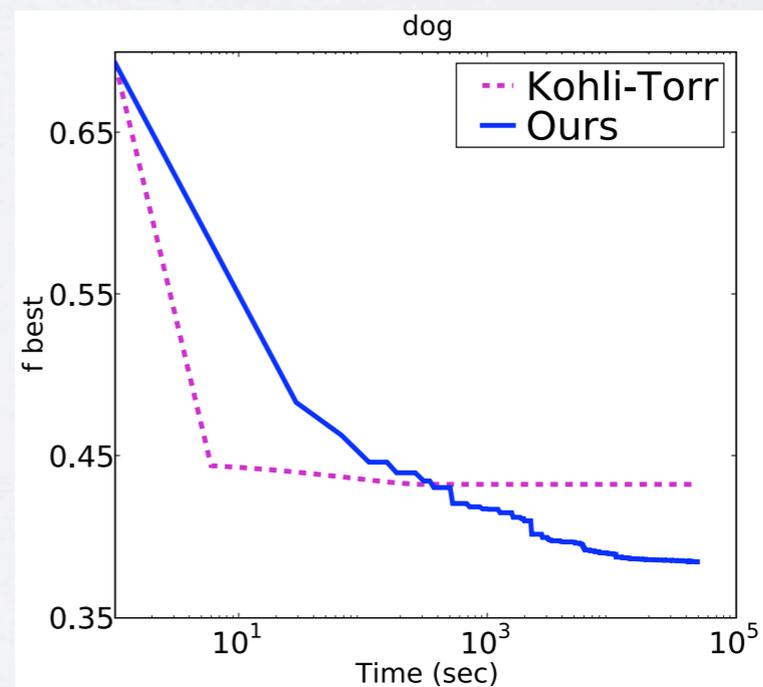
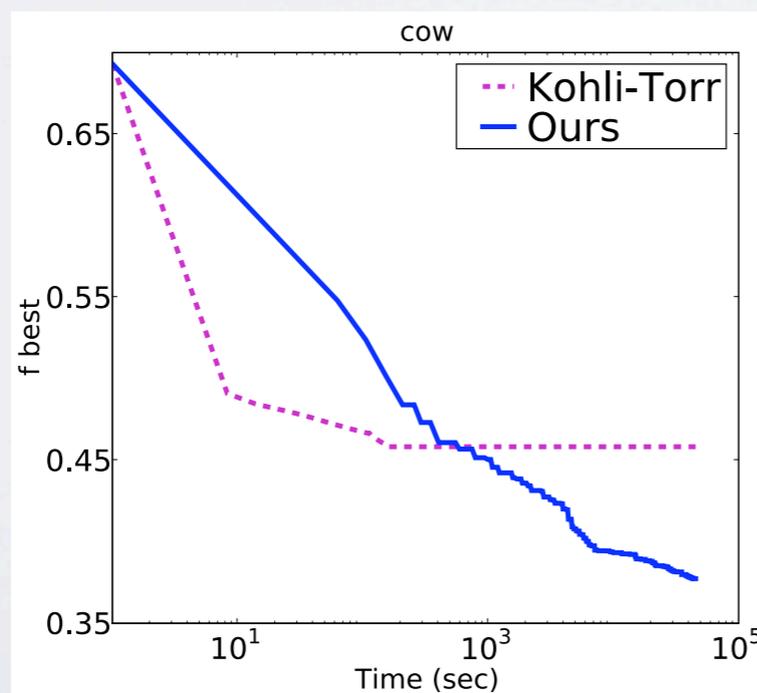
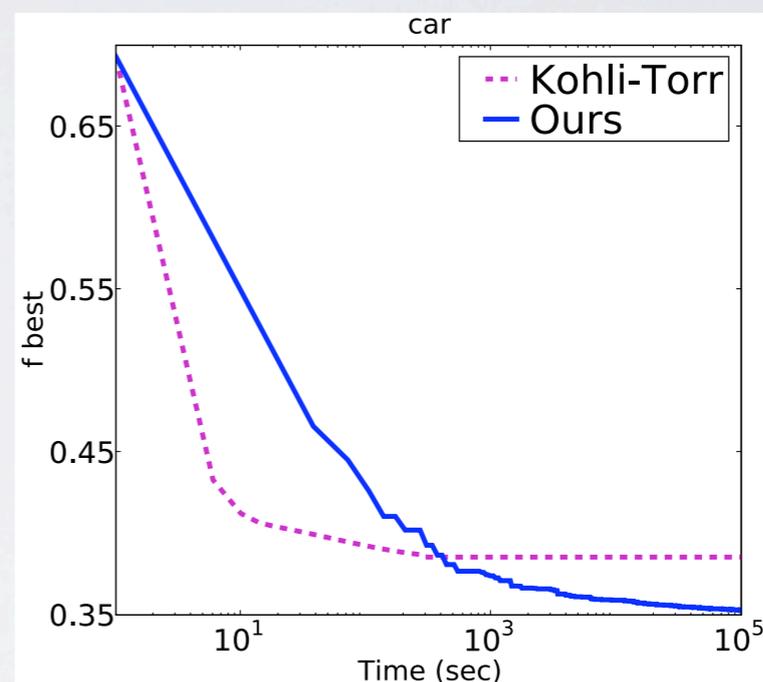
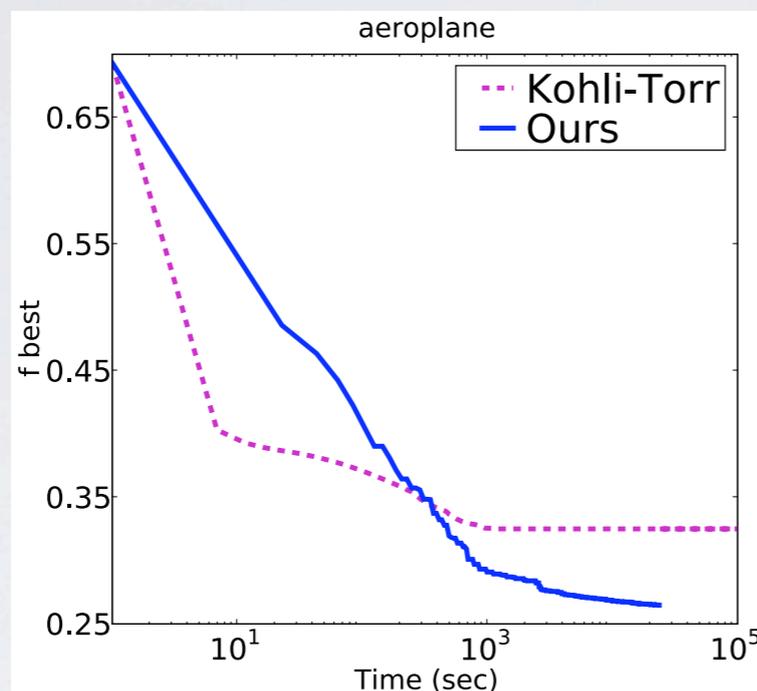
Do we get better training likelihoods by optimizing the model with the min-marginal objective?

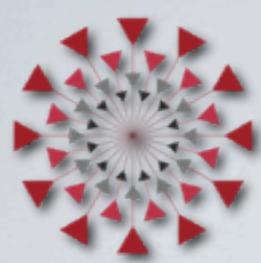




# EVALUATING TRAINING PROCEDURES

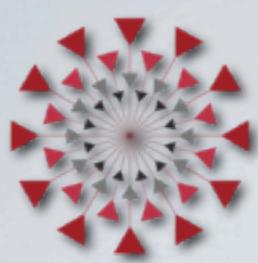
Does training time improve using the gradient of the log likelihood vs procedure of K&T?





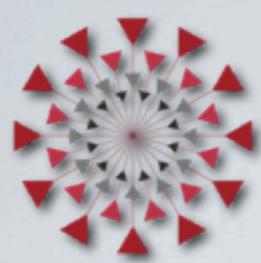
# EVALUATING TEST-TIME PERFORMANCE

		<b>Log Lik</b>	<b>Accuracy</b>	<b>AUC</b>
<b>Aero</b>	KT-final	-.35 (-.29)	87.3 (89.7)	.85 (.87)
	KT- $f_{best}$	-.32 (-.28)	88.1 (89.9)	.85 (.87)
	Ours	<b>-.26 (-.24)</b>	<b>88.9 (90.4)</b>	<b>.90 (.90)</b>
<b>Car</b>	KT-final	-.48 (-.65)	84.1 (78.7)	.69 (.65)
	KT- $f_{best}$	-.39 (-.51)	<b>86.2 (80.0)</b>	.66 (.62)
	Ours	<b>-.35 (-.51)</b>	86.1 ( <b>81.4</b> )	<b>.76 (.65)</b>
<b>Cow</b>	KT-final	-.54 (-.64)	79.7 (75.9)	.76 (.77)
	KT- $f_{best}$	-.47 (-.52)	80.9 (76.7)	.66 (.65)
	Ours	<b>-.38 (-.41)</b>	<b>82.5 (79.9)</b>	<b>.84 (.82)</b>
<b>Dog</b>	KT-final	-.52 (-.45)	81.8 (84.7)	.64 (.66)
	KT- $f_{best}$	-.43 (-.38)	<b>84.1 (86.7)</b>	.62 (.66)
	Ours	<b>-.38 (-.34)</b>	84.0 ( <b>86.8</b> )	<b>.76 (.79)</b>



# EVALUATING TEST-TIME PERFORMANCE

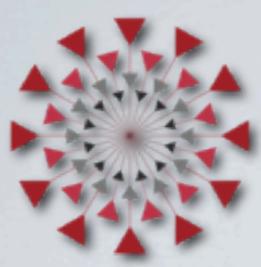
<b>Image</b>			
<b>True Label</b>			
<b>KT-final</b>			
<b>KT-<math>f_{best}</math></b>			
<b>Ours</b>			



# INCORPORATING TEST-TIME LOSSES

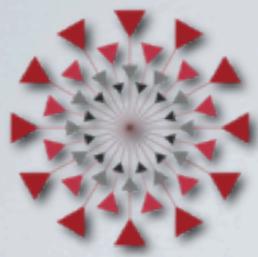
If we must produce **one** test-time segmentation, better-calibrated probabilities enable us to choose the segmentation that **minimizes expected loss**.

		$n/U$ Before	$n/U$ After	Change
<b>Aero</b>	KT-final	63.4	63.9	.5
	KT- $f_{best}$	60.7	61.8	1.1
	Ours	<b>64.1</b>	<b>66.6</b>	<b>2.5</b>
<b>Car</b>	KT-final	<b>43.5</b>	44.2	.7
	KT- $f_{best}$	40.7	43.3	2.6
	Ours	41.8	<b>45.9</b>	<b>4.1</b>
<b>Cow</b>	KT-final	<b>51.2</b>	51.7	.5
	KT- $f_{best}$	40.7	47.1	<b>6.4</b>
	Ours	46.5	<b>52.8</b>	6.3
<b>Dog</b>	KT-final	<b>52.1</b>	52.0	-.1
	KT- $f_{best}$	44.5	47.2	2.7
	Ours	48.9	<b>55.3</b>	<b>6.4</b>



# INCORPORATING TEST-TIME LOSSES

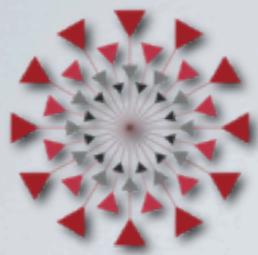
<b>Image</b>			
<b>True Label</b>			
<b>Ours</b>			
<b>Before</b>			
<b>After</b>			



# SUMMARY

---

- ▶ MRF doesn't stand for "Magical Random Field"
- ▶ Efficient optimization procedures can provide a useful tool in the generative modeling toolbox for structured prediction.
- ▶ Train your models to answer the questions you'll ask them at test time!
- ▶ If the full representational power of the model will never be used, there may be a significant computational gain in exactly fitting a model that is "just big enough".



# THANKS

---

- ▶ Danny Tarlow (University of Toronto)
- ▶ Rich Zemel (University of Toronto)



- ▶ D. Tarlow, R.P. Adams, R.S. Zemel. Randomized Optimum Models for Structured Prediction. AISTATS 2012.
- ▶ D. Tarlow, R.P. Adams. Revisiting Uncertainty in Graph Cut Solutions. CVPR 2012.