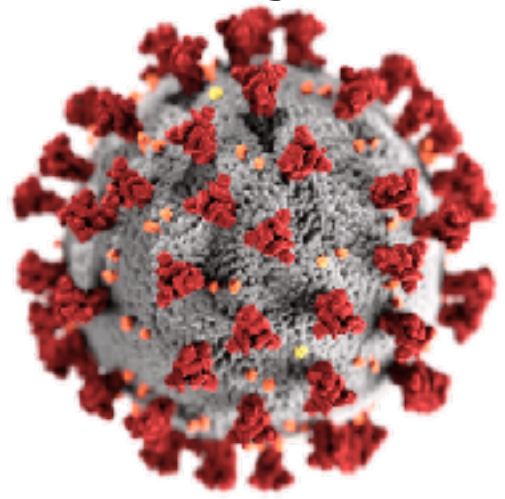


Variational Bayes and beyond: Foundations of scalable Bayesian inference

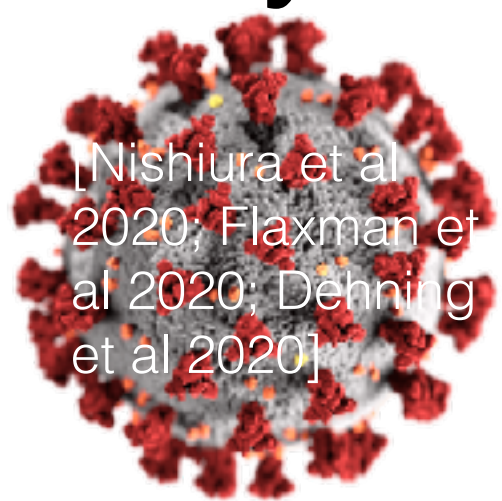
Tamara Broderick
Associate Professor
MIT

Bayesian inference

Bayesian inference

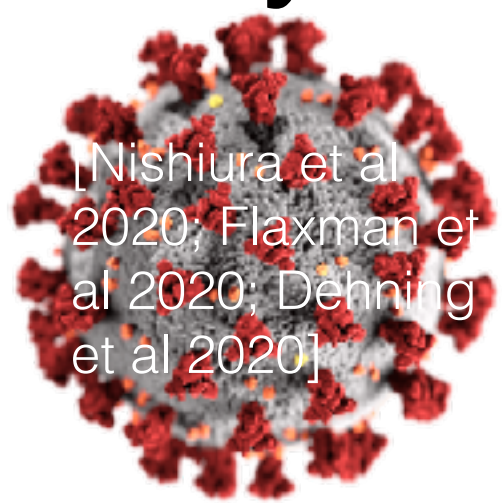


Bayesian inference

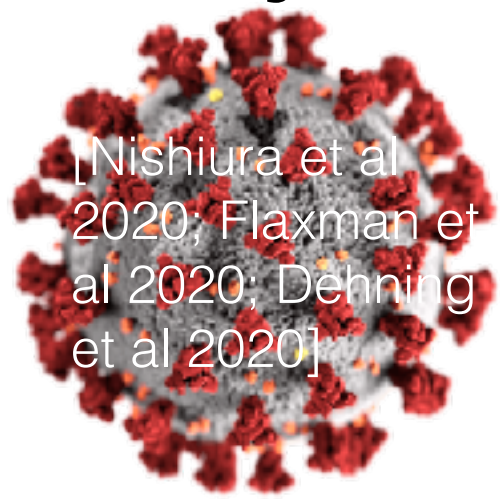


[Nishiura et al
2020; Flaxman et
al 2020; Dehning
et al 2020]

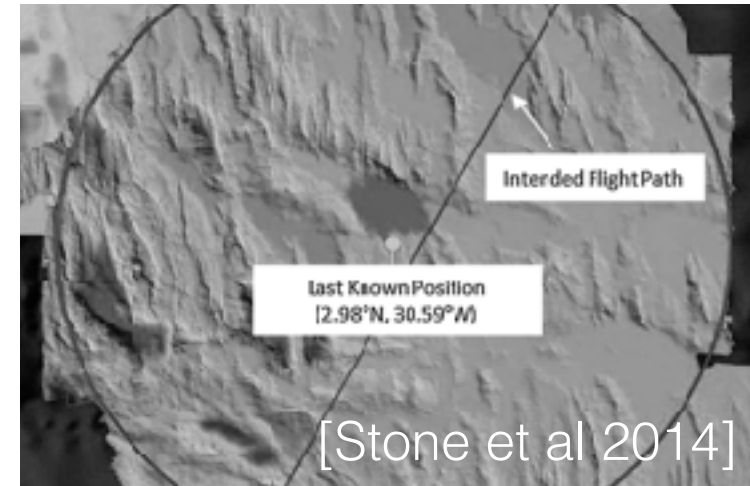
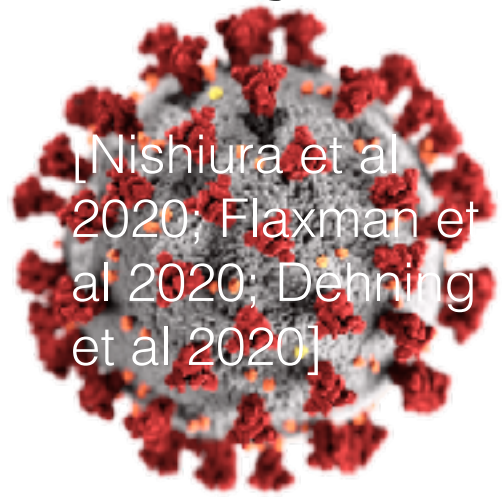
Bayesian inference



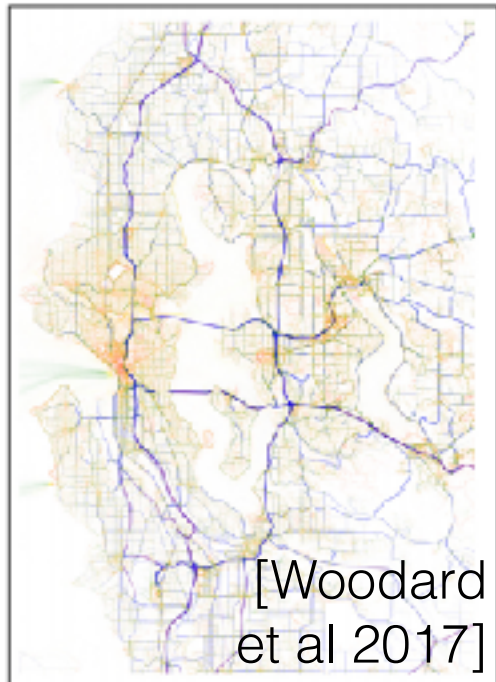
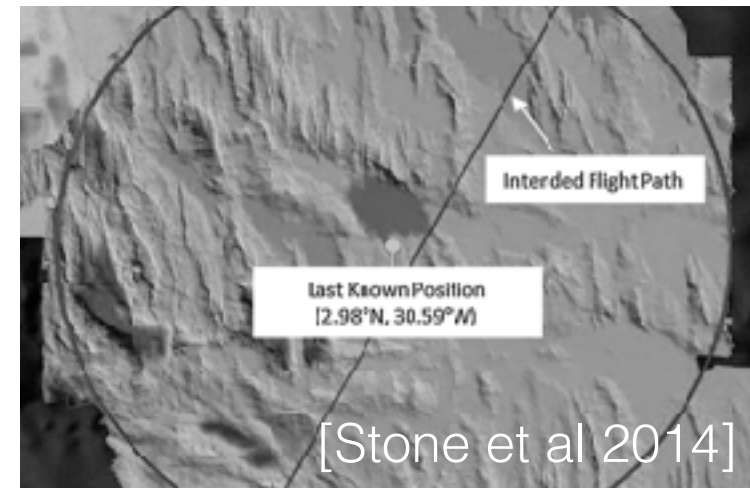
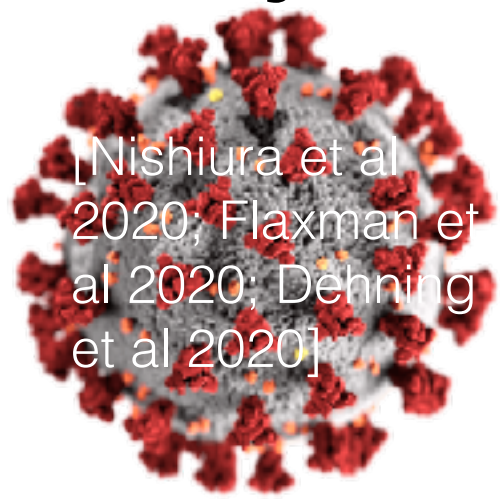
Bayesian inference



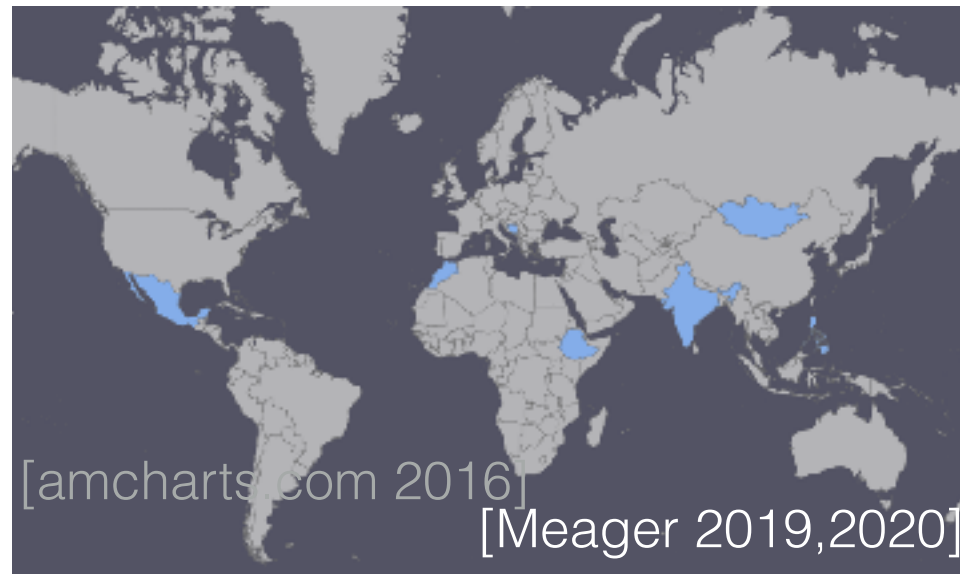
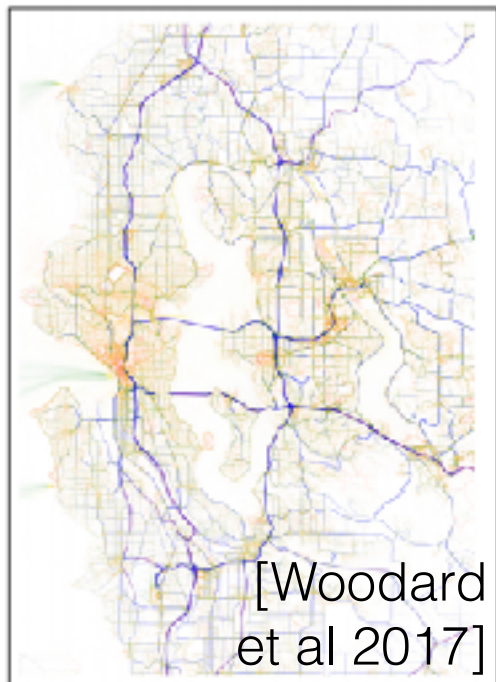
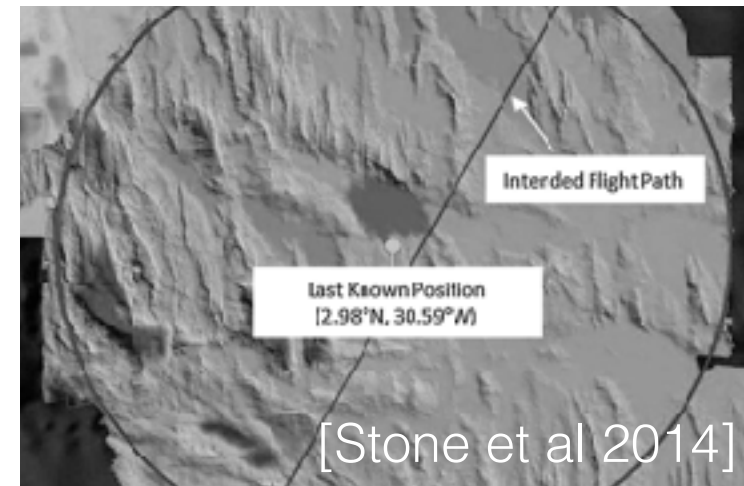
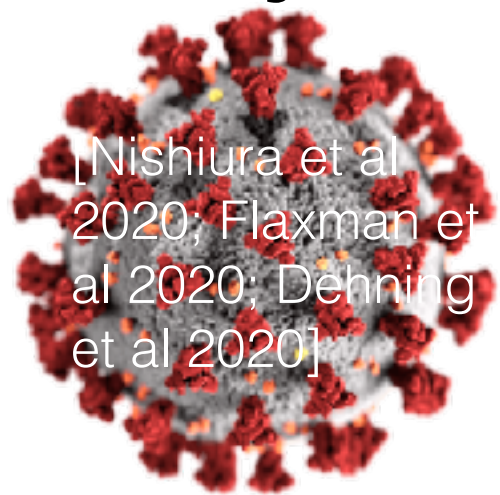
Bayesian inference



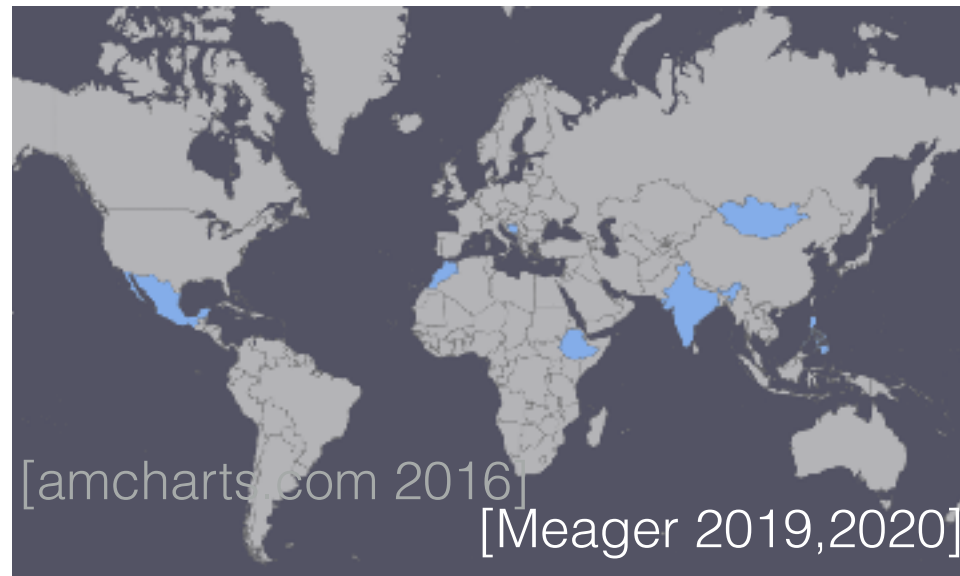
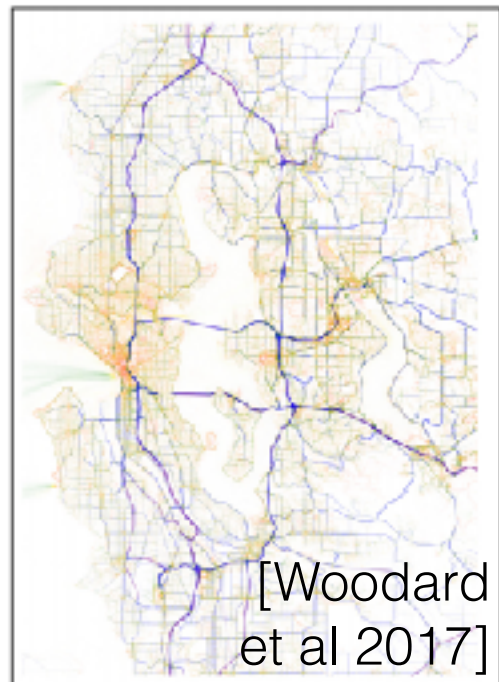
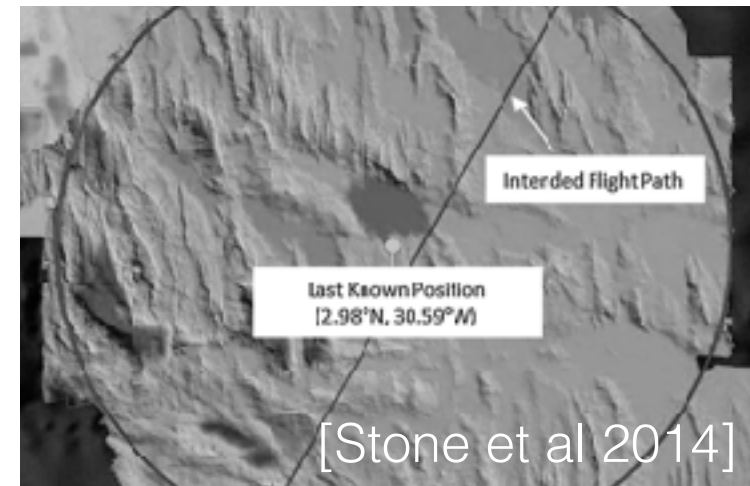
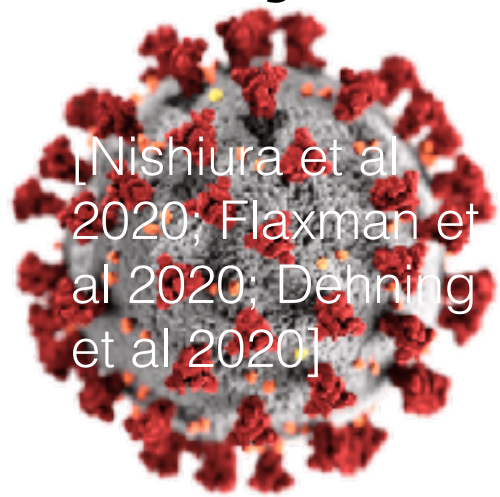
Bayesian inference



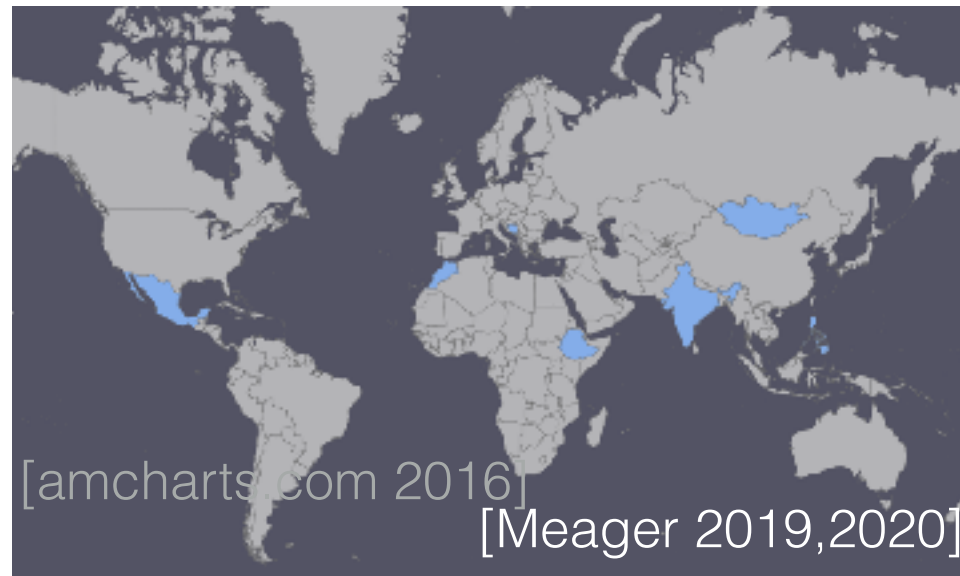
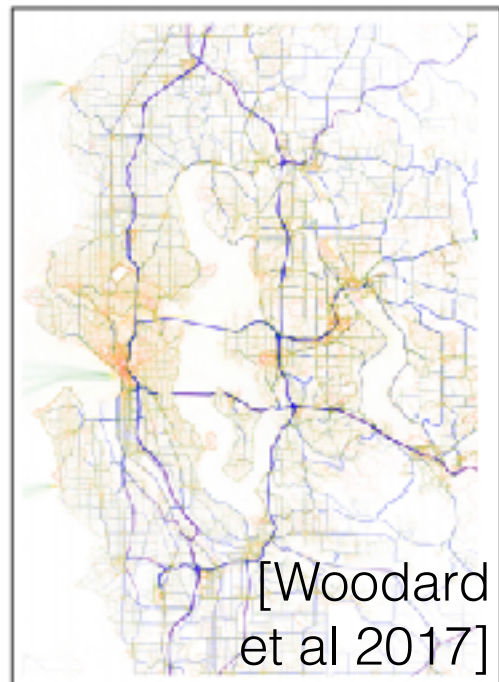
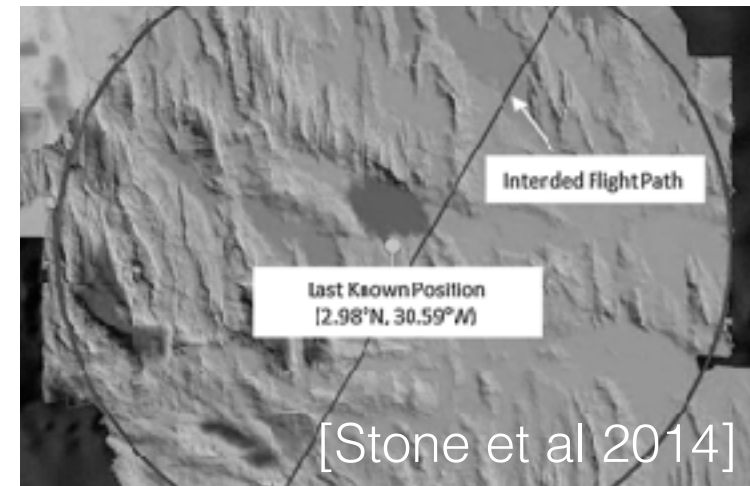
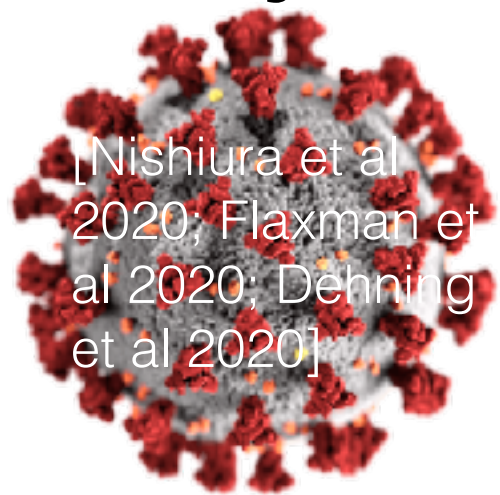
Bayesian inference



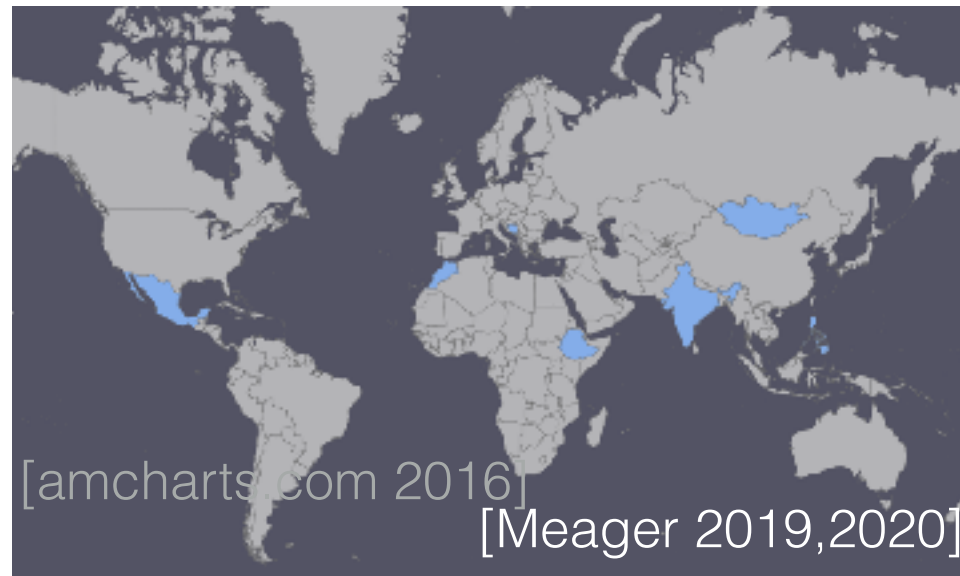
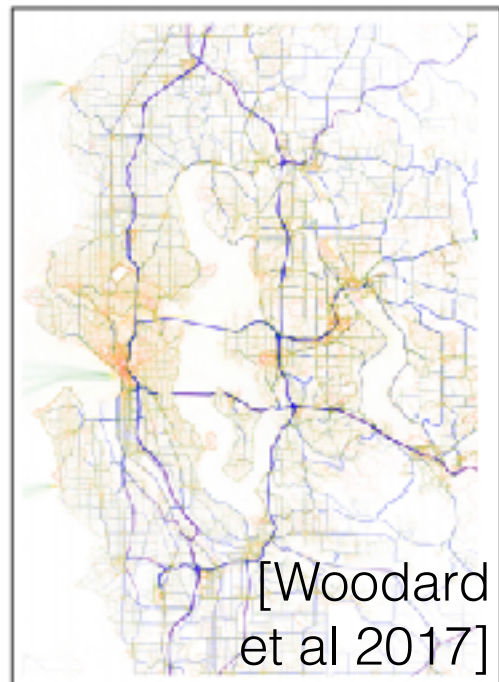
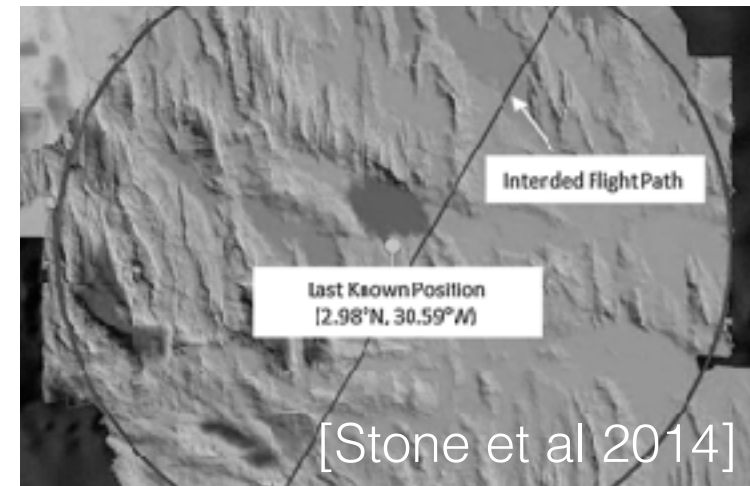
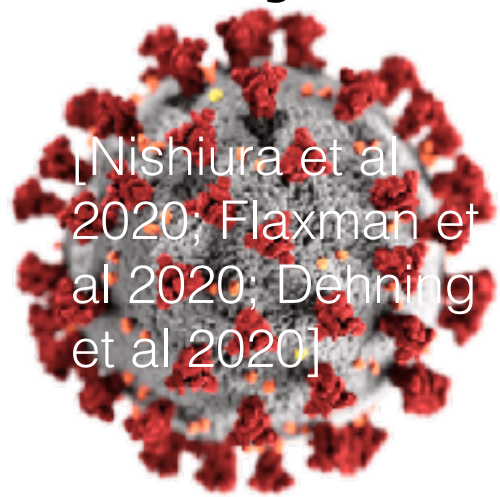
Bayesian inference



Bayesian inference

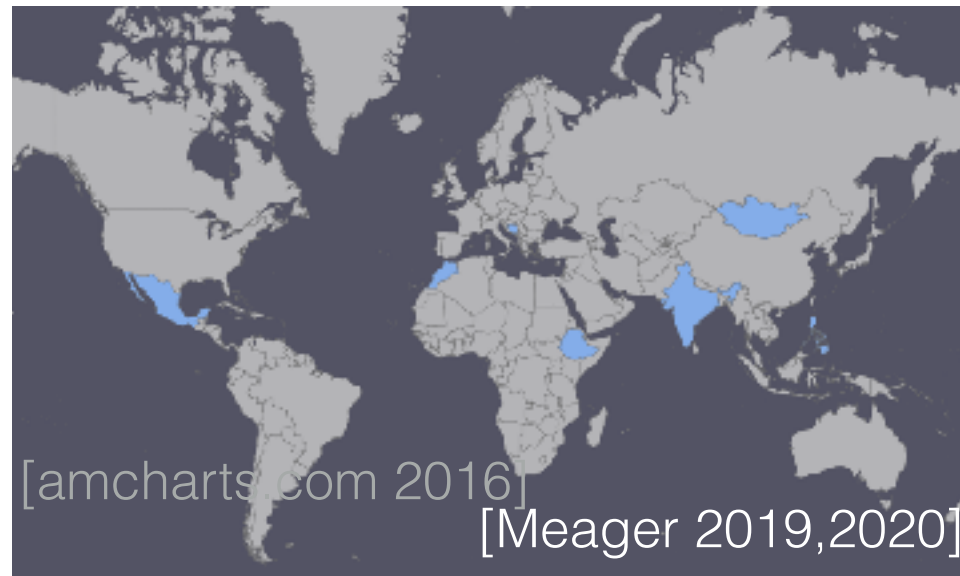
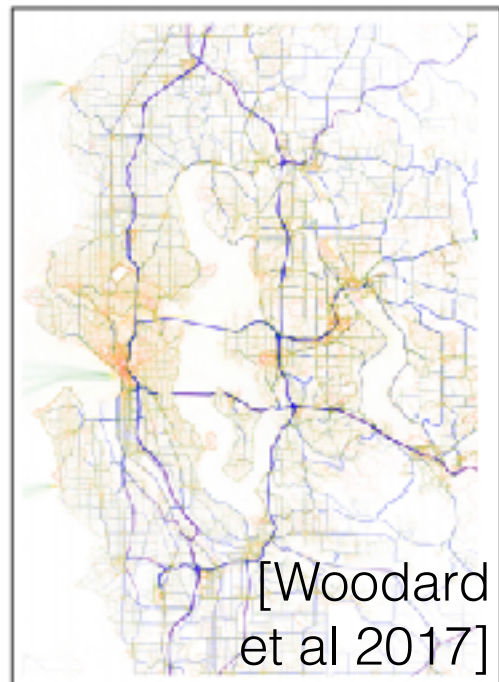
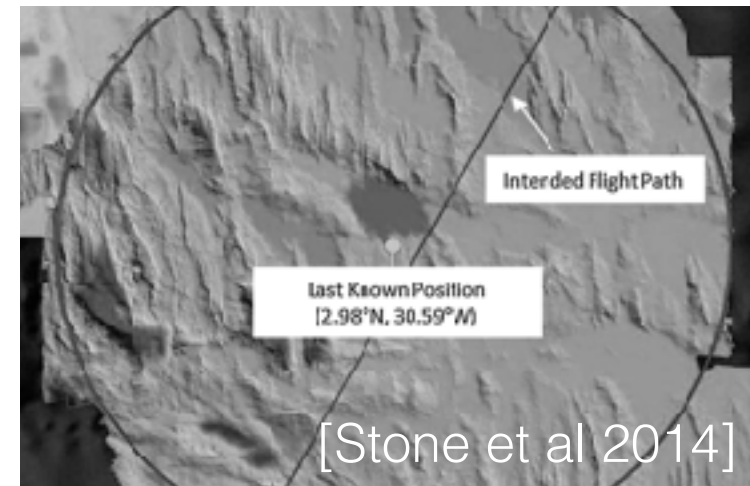
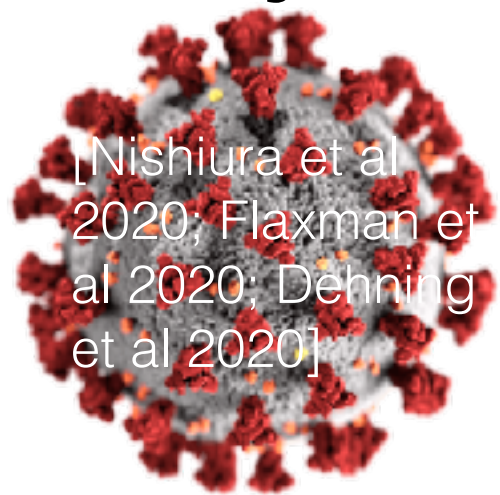


Bayesian inference



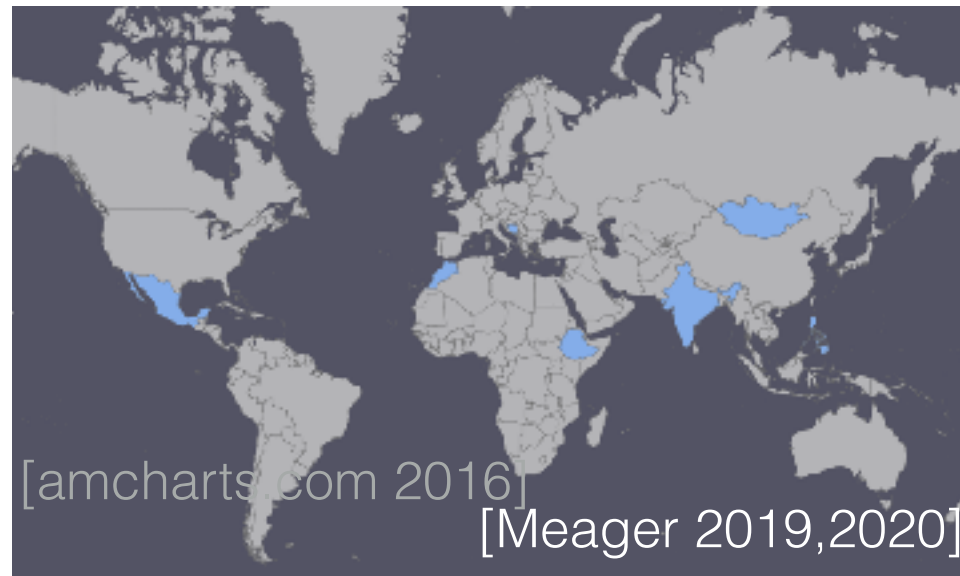
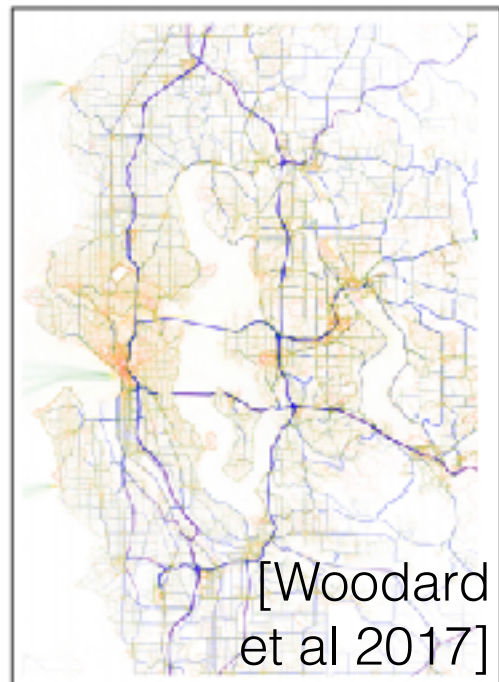
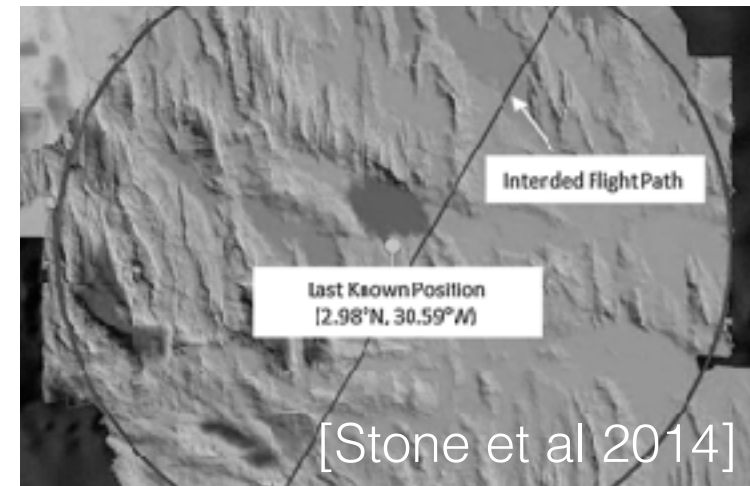
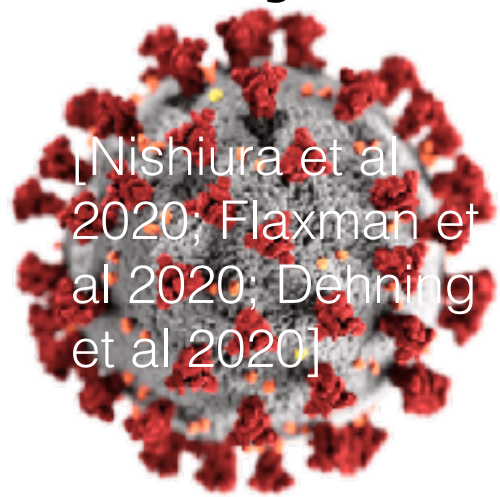
- Goals: good point estimates, uncertainty estimates

Bayesian inference



- Goals: good point estimates, uncertainty estimates
 - More: interpretable, modular, expert info

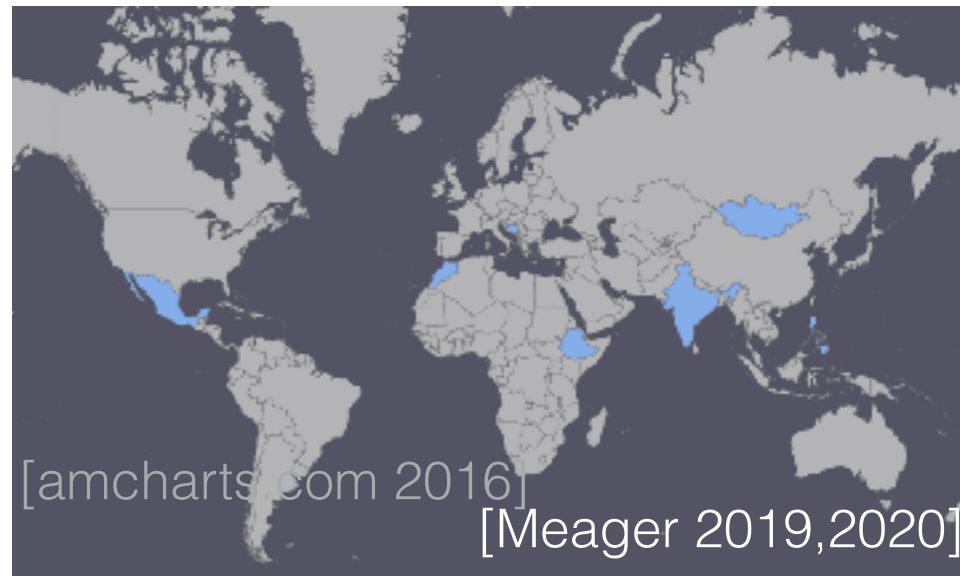
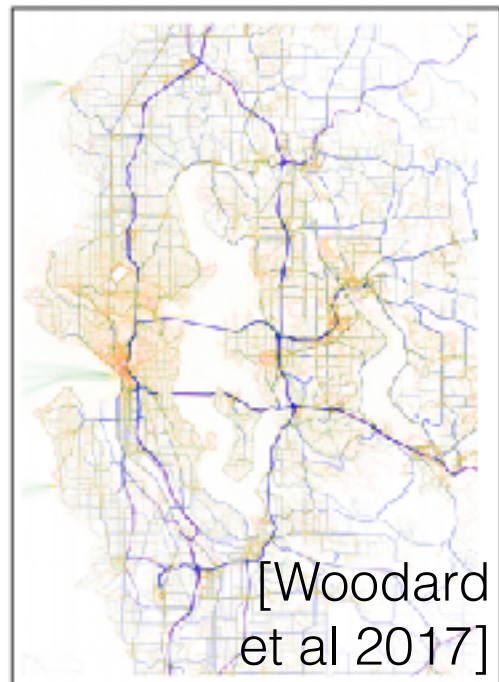
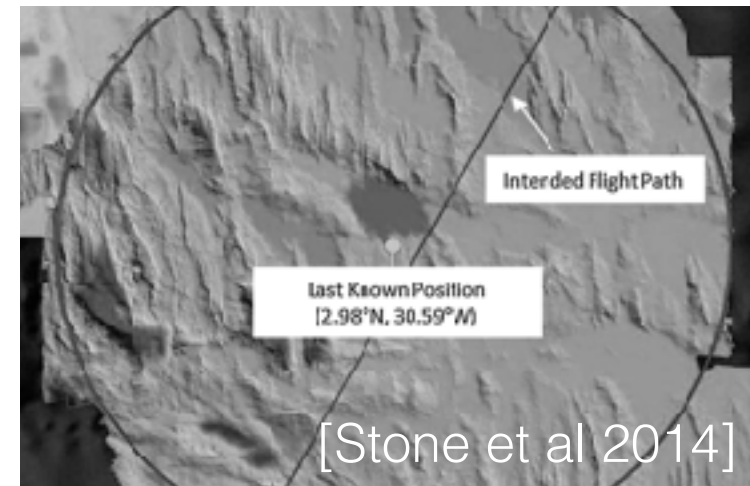
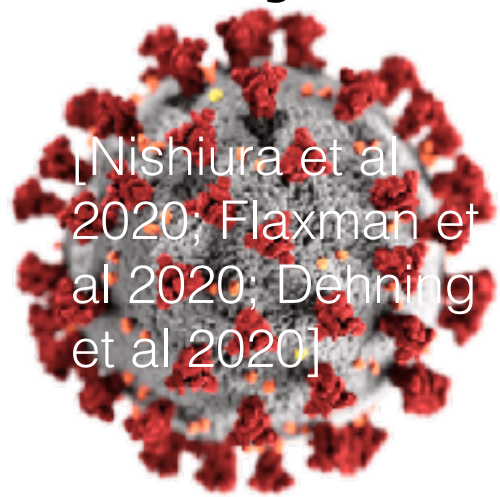
Bayesian inference



[mc-stan.org]

- Goals: good point estimates, uncertainty estimates
- More: interpretable, modular, expert info

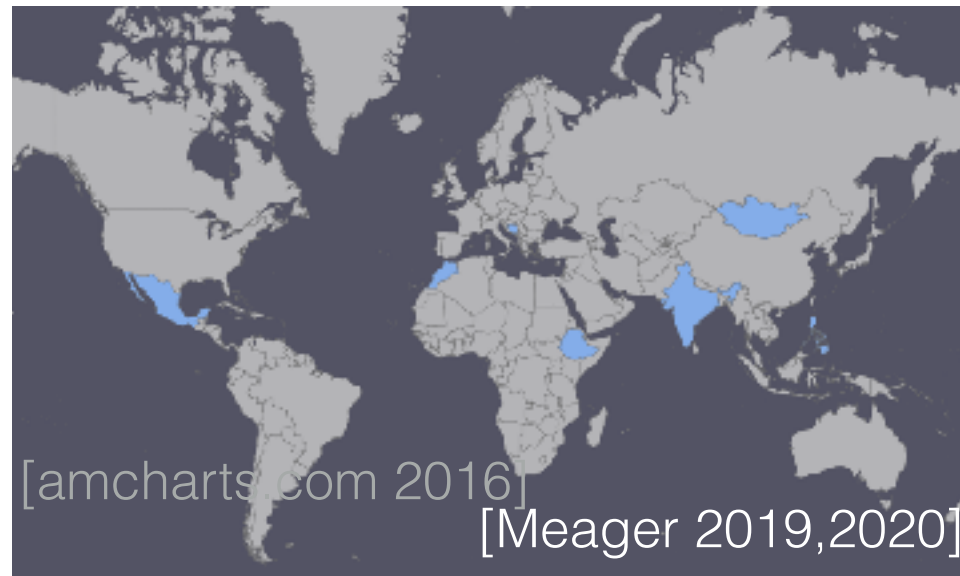
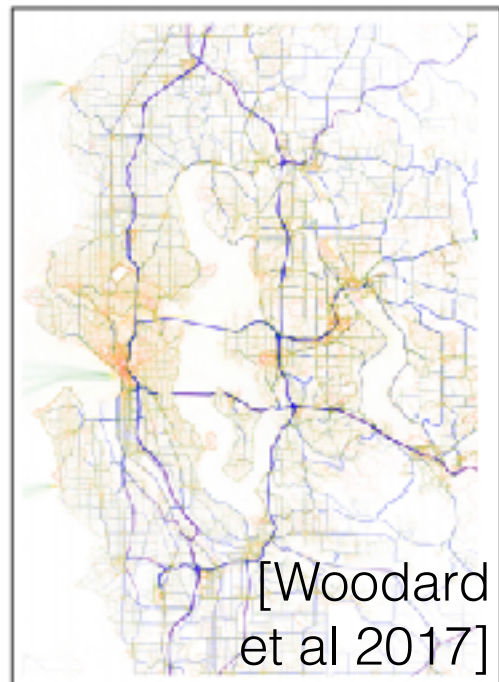
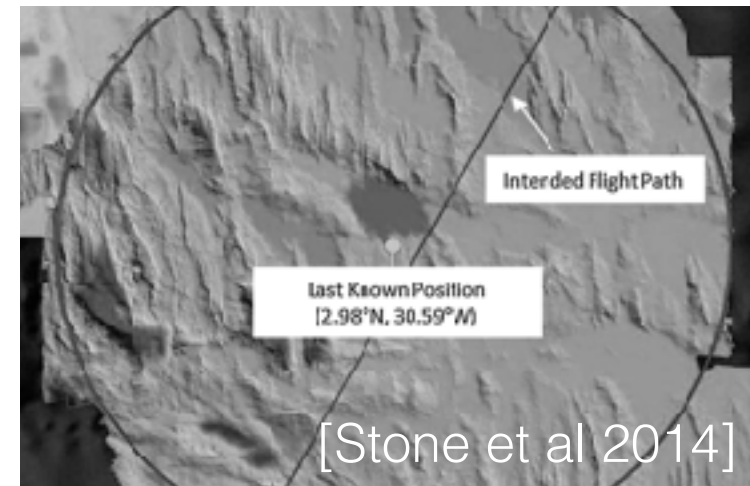
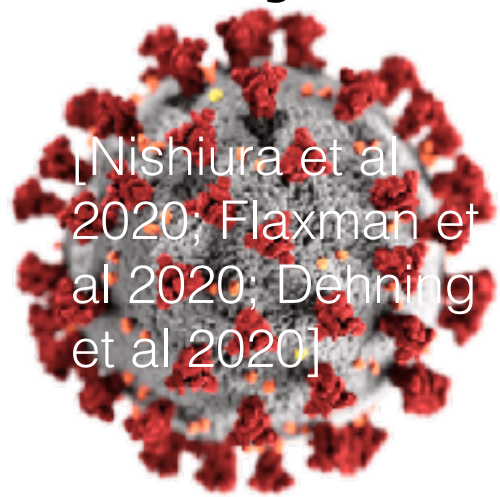
Bayesian inference



[mc-stan.org]

- Goals: good point estimates, uncertainty estimates
 - More: interpretable, modular, expert info
- Challenge: speed (compute, user), reliable inference

Bayesian inference



[mc-stan.org]

- Goals: good point estimates, uncertainty estimates
 - More: interpretable, modular, expert info
- Challenge: speed (compute, user), reliable inference
- Uncertainty doesn't have to disappear in large data sets

Variational Bayes

Variational Bayes

- Modern problems: often large data, large dimensions

Variational Bayes

- Modern problems: often large data, large dimensions
- Variational Bayes can be very fast

Variational Bayes

- Modern problems: often large data, large dimensions
- Variational Bayes can be very fast

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

[Blei et al
2003]

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

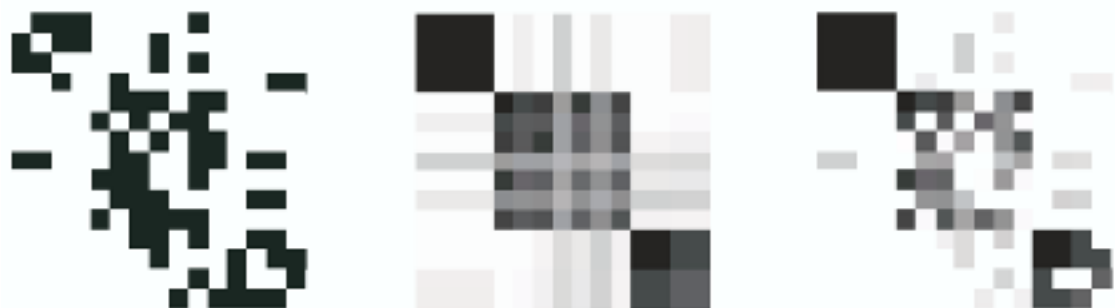
Variational Bayes

- Modern problems: often large data, large dimensions
- Variational Bayes can be very fast

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

[Blei et al
2003]

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



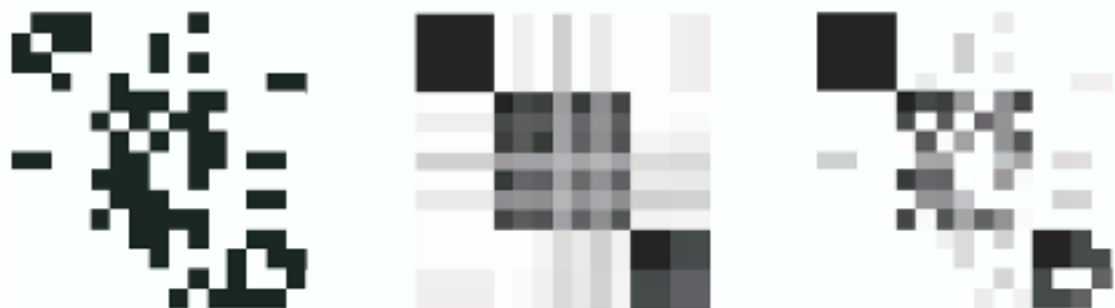
Variational Bayes

- Modern problems: often large data, large dimensions
- Variational Bayes can be very fast

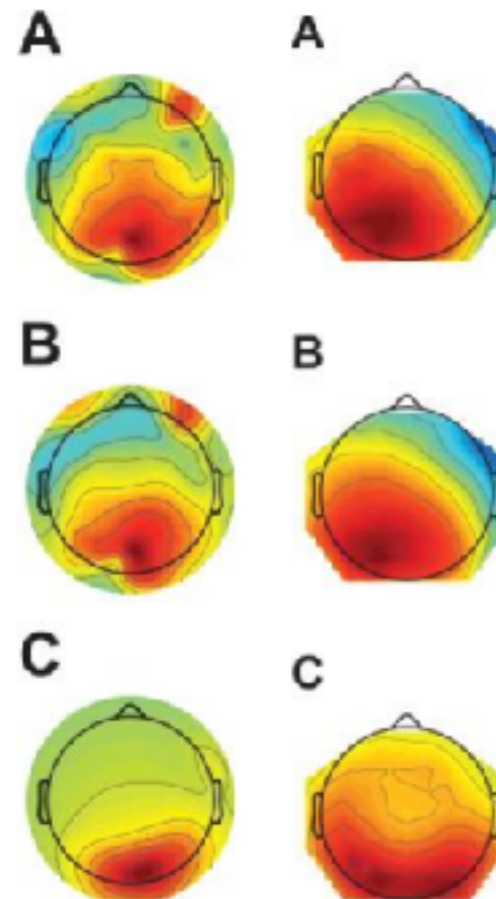
"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

[Blei et al
2003]

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



[Airola et al 2008]



[Gershman et al 2014]

[Blei et al 2018]

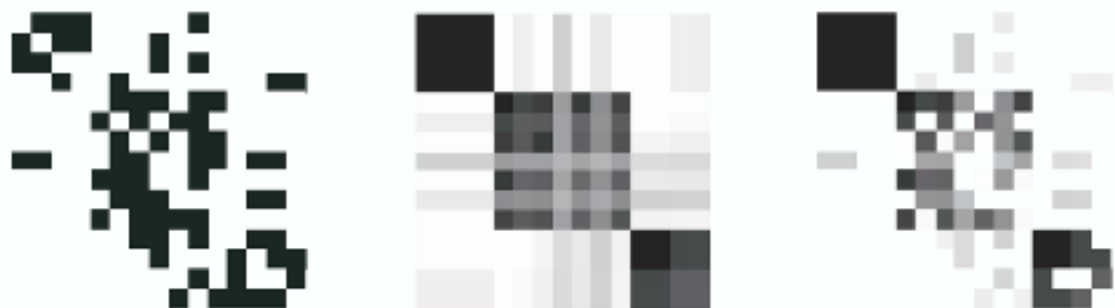
Variational Bayes

- Modern problems: often large data, large dimensions
- Variational Bayes can be very fast

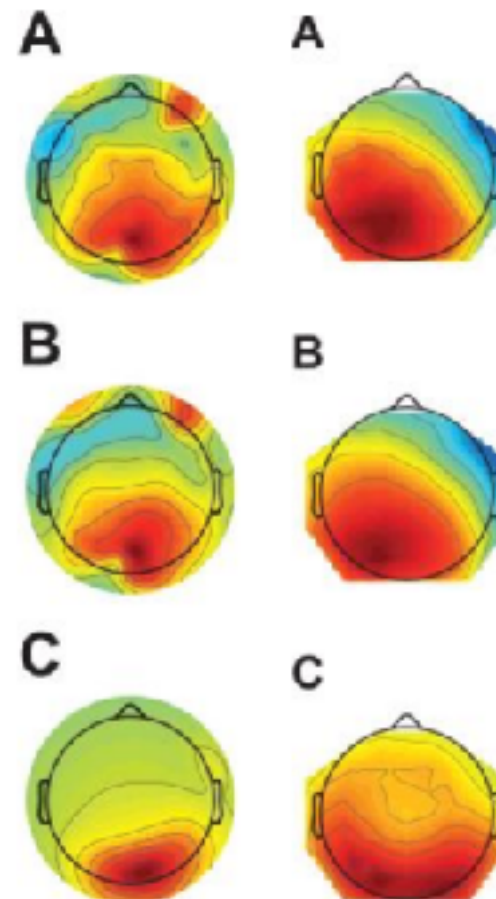
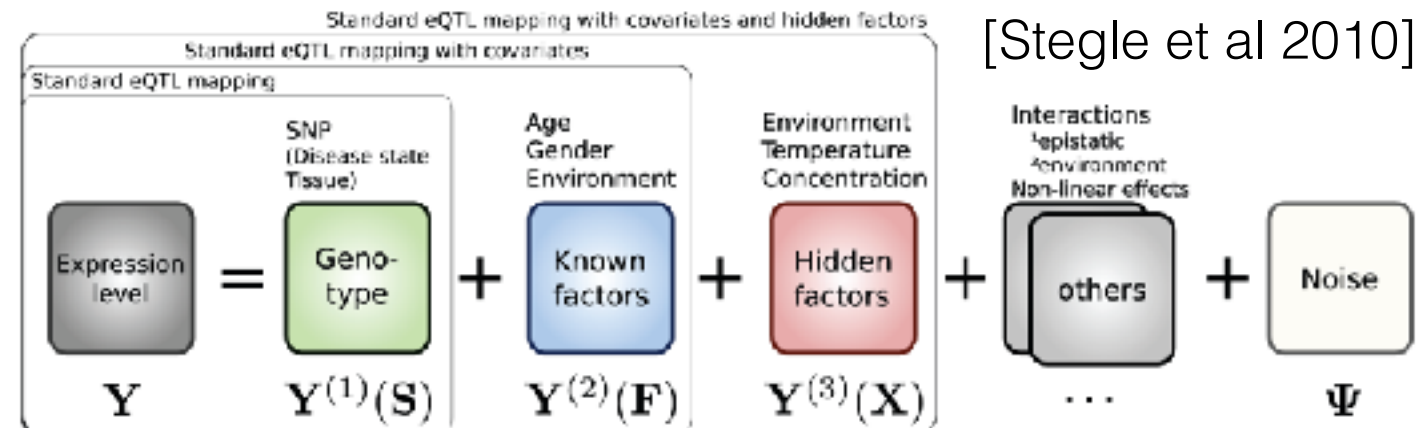
"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

[Blei et al 2003]

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



[Airoldi et al 2008]



[Gershman et al 2014]

[Blei et al 2018]

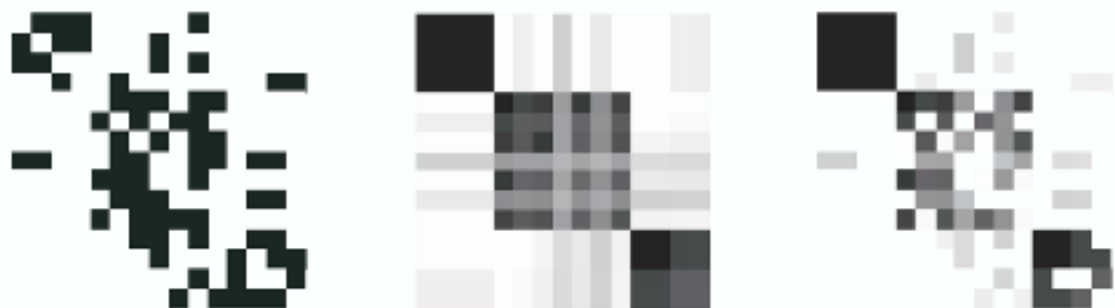
Variational Bayes

- Modern problems: often large data, large dimensions
- Variational Bayes can be very fast

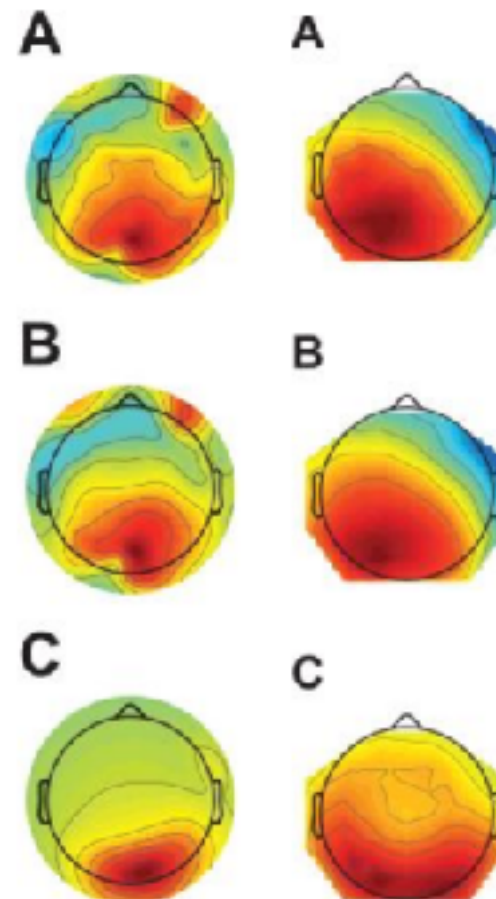
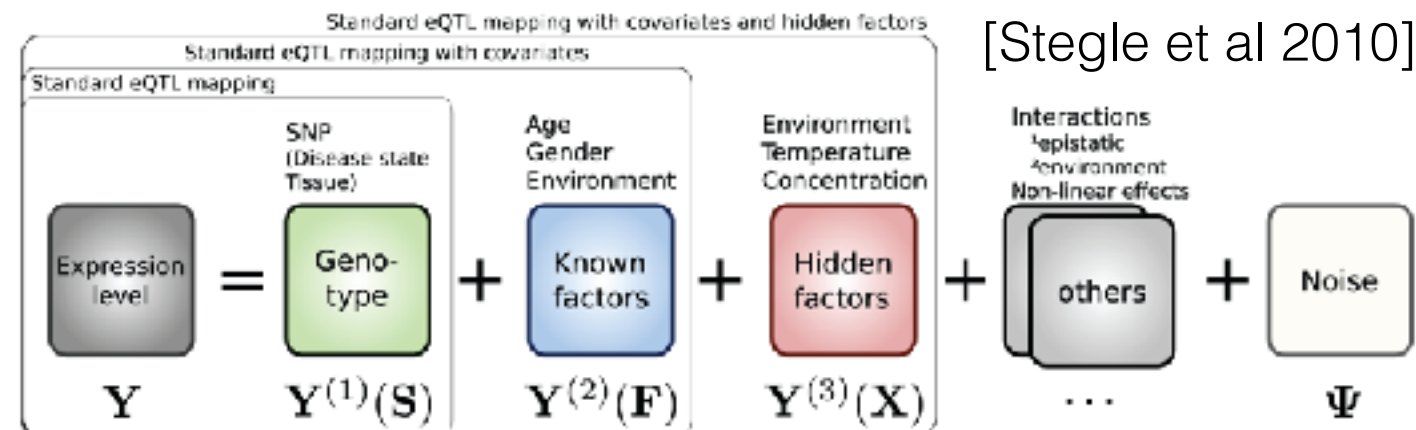
"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

[Blei et al 2003]

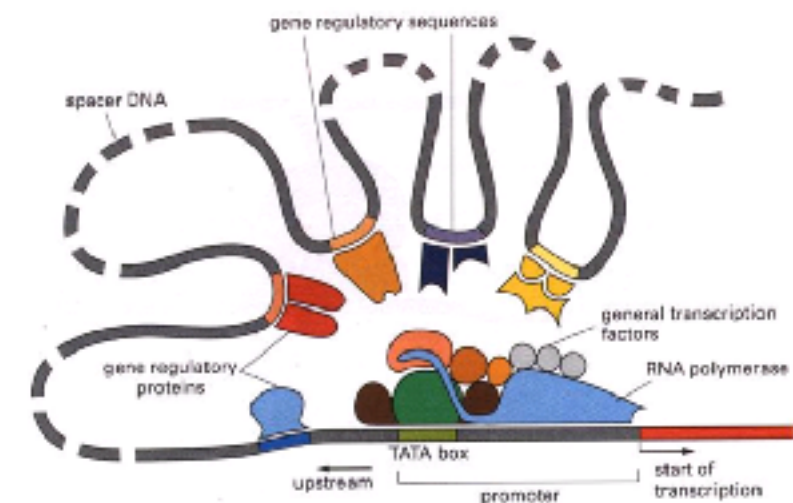
The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



[Airoldi et al 2008]



[Gershman et al 2014]



[Xing et al 2004]

[Xing 2003]

[Blei et al 2018]

Roadmap

- Bayes & Approximate Bayes review

Roadmap

- Bayes & Approximate Bayes review
- What is:
 - Variational Bayes (VB)
 - Mean-field variational Bayes (MFVB)

Roadmap

- Bayes & Approximate Bayes review
- What is:
 - Variational Bayes (VB)
 - Mean-field variational Bayes (MFVB)
- Why use VB?

Roadmap

- Bayes & Approximate Bayes review
- What is:
 - Variational Bayes (VB)
 - Mean-field variational Bayes (MFVB)
- Why use VB?
- When can we trust VB?

Roadmap

- Bayes & Approximate Bayes review
- What is:
 - Variational Bayes (VB)
 - Mean-field variational Bayes (MFVB)
- Why use VB?
- When can we trust VB?
- Where do we go from here?

Roadmap

- Bayes & Approximate Bayes review
- What is:
 - Variational Bayes (VB)
 - Mean-field variational Bayes (MFVB)
- Why use VB?
- When can we trust VB?
- Where do we go from here?

Bayesian inference

Bayesian inference

parameters

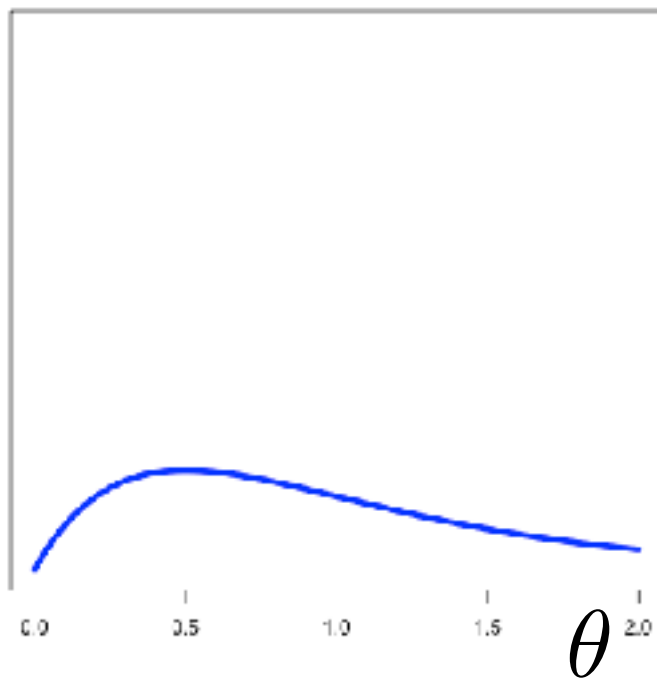
 θ

Bayesian inference

parameters
↓
 $p(\theta)$
prior

Bayesian inference

parameters
↓
 $p(\theta)$
prior



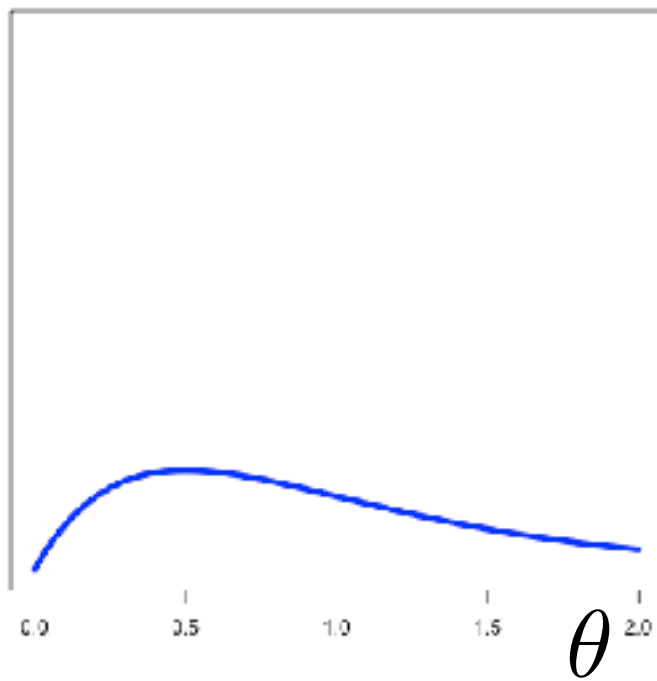
Bayesian inference

parameters



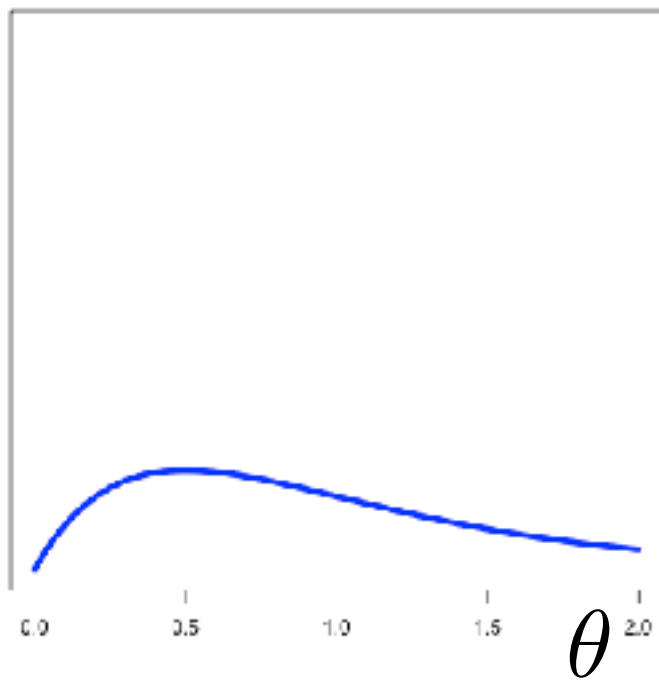
$$p(y_{1:N}|\theta)p(\theta)$$

likelihood prior



Bayesian inference

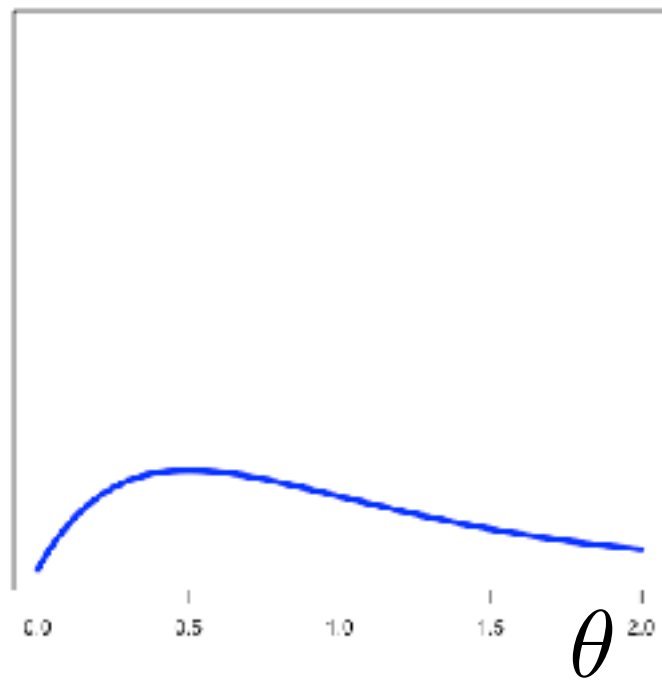
data parameters
 $p(y_{1:N}|\theta)p(\theta)$
likelihood prior



Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior



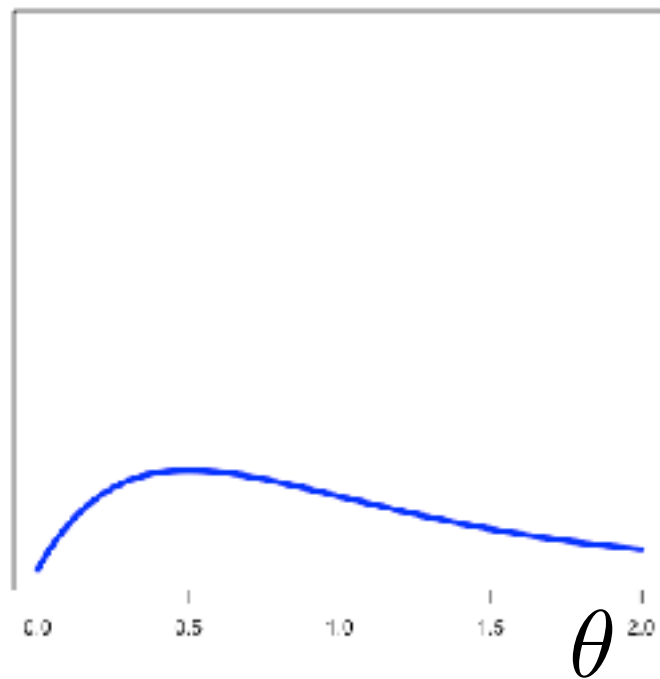
Bayesian inference


data

parameters

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior

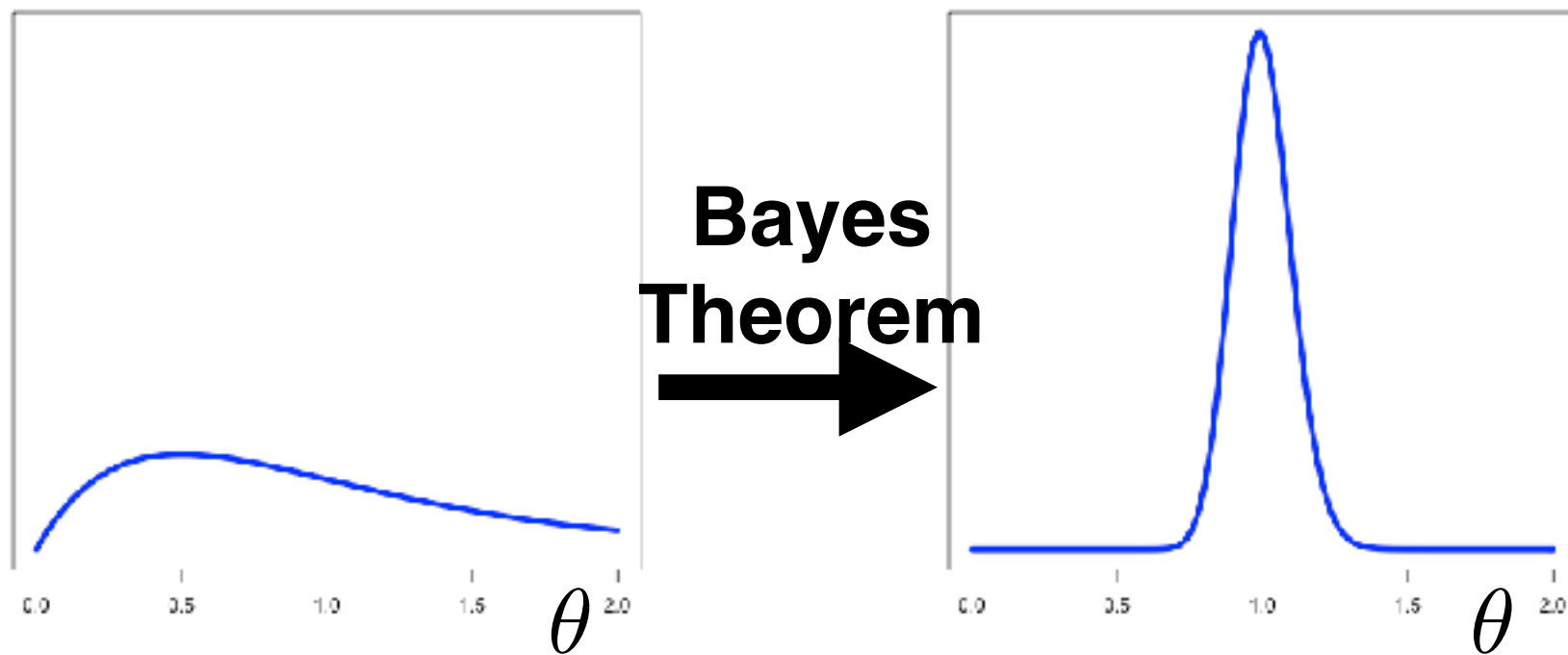


**Bayes
Theorem** 

Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

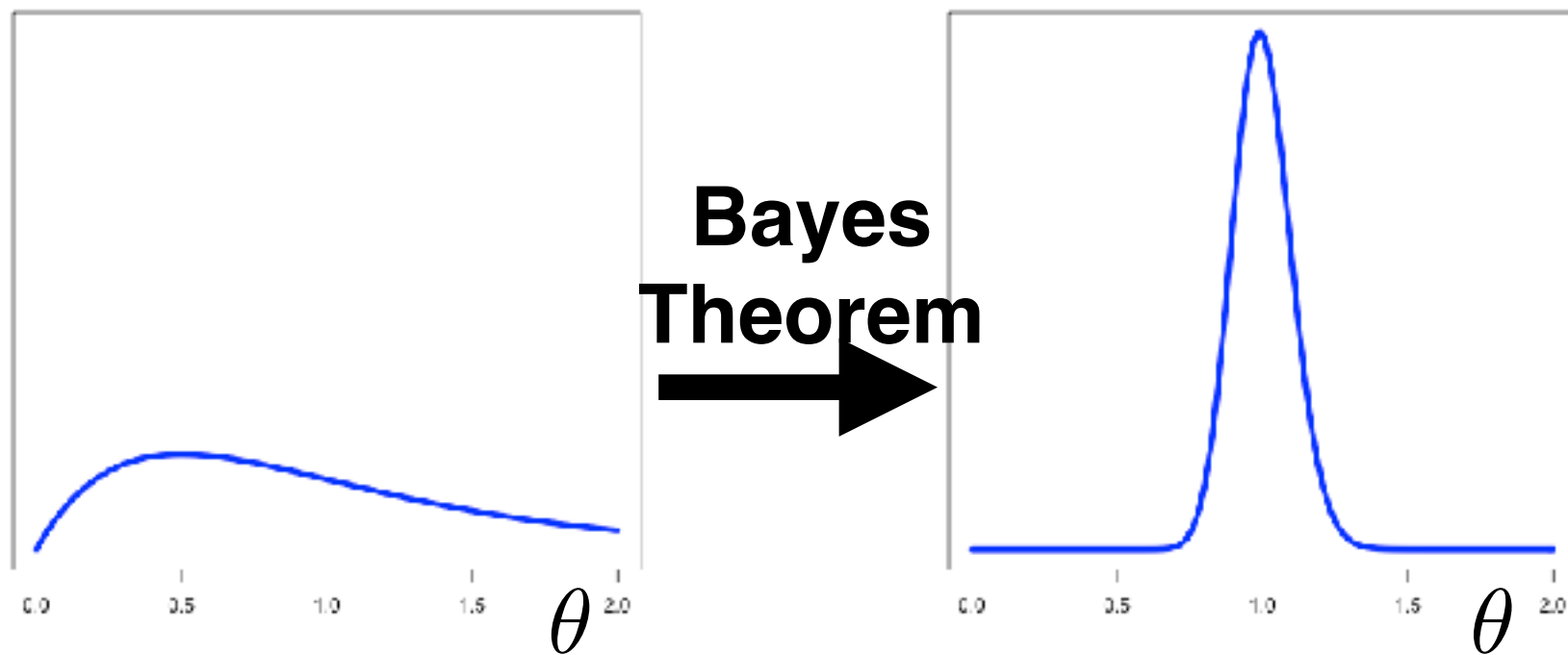
posterior likelihood prior



Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior

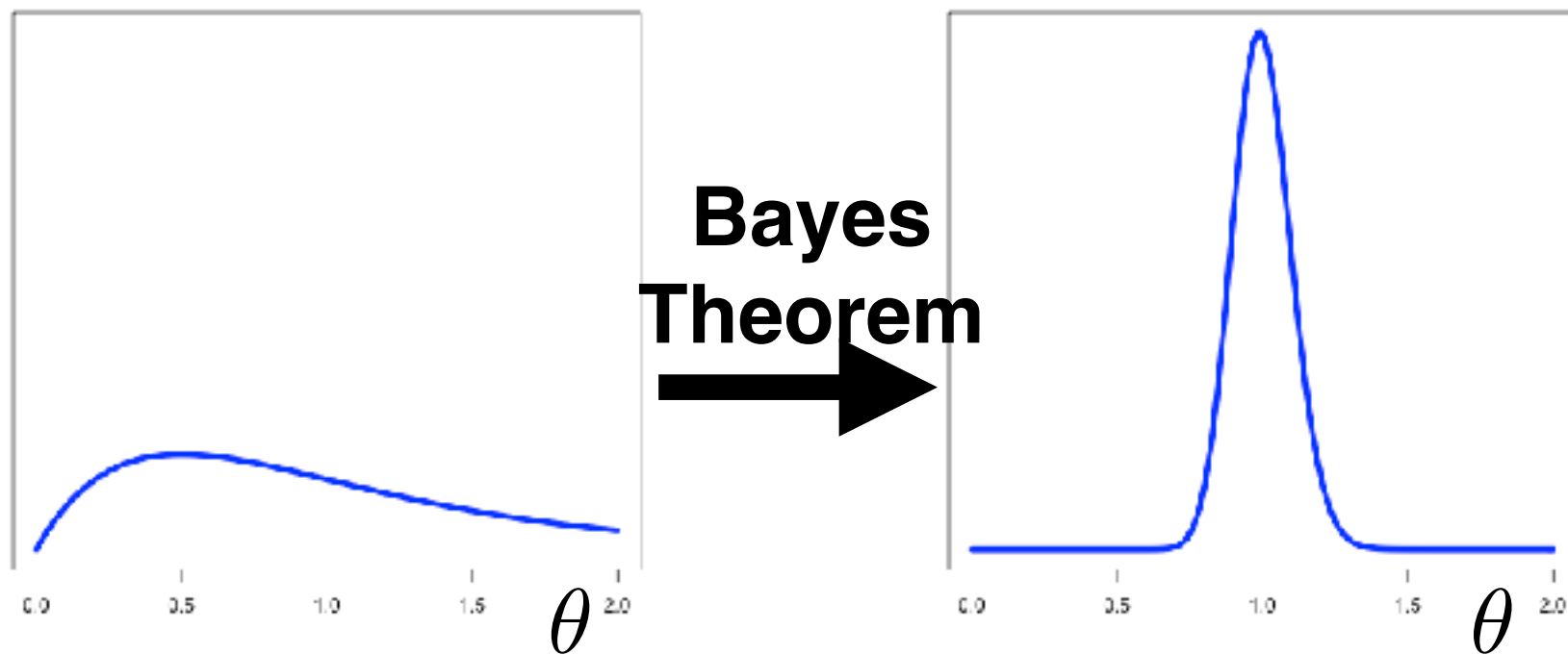


1. Build a model: choose prior & choose likelihood

Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior

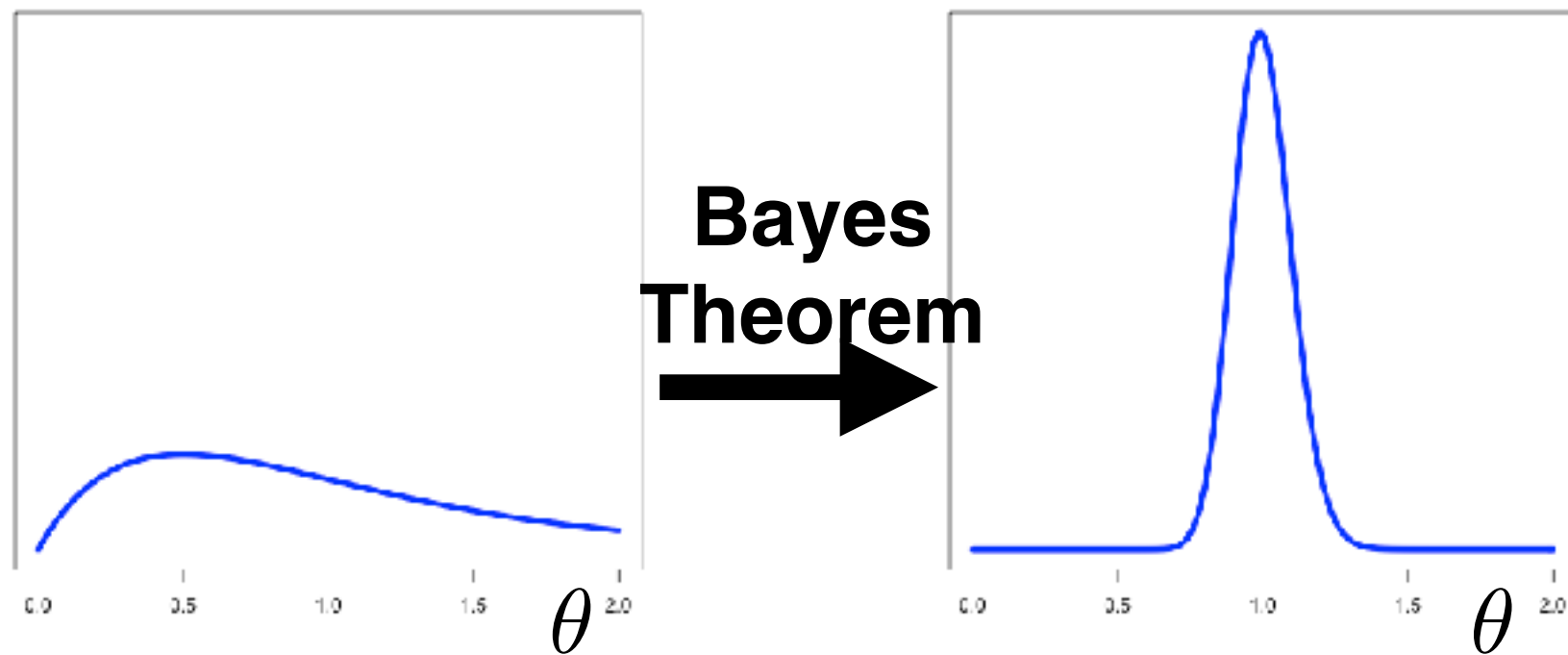


1. Build a model: choose prior & choose likelihood
2. Compute the posterior

Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior

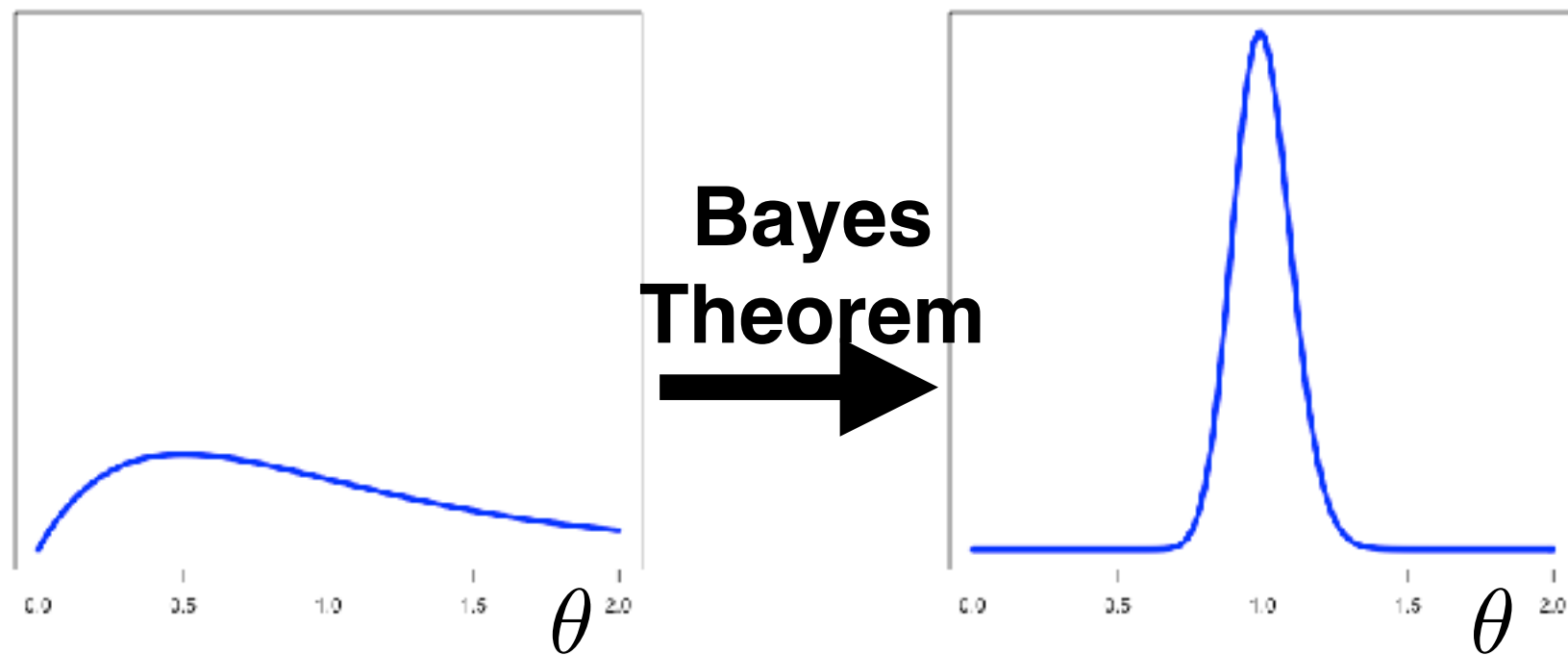


1. Build a model: choose prior & choose likelihood
2. Compute the posterior
3. Report a summary, e.g. posterior means and (co)variances

Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior



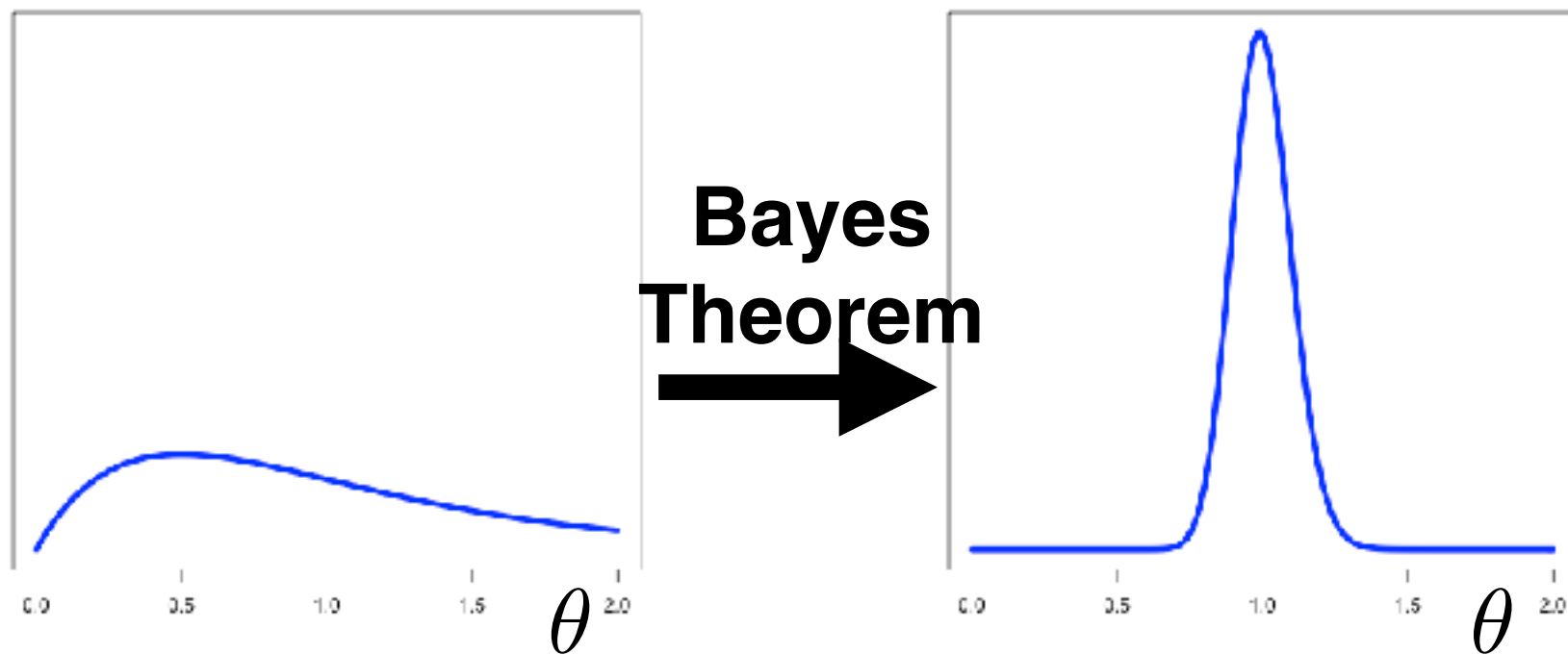
1. Build a model: choose prior & choose likelihood
 2. Compute the posterior
 3. Report a summary, e.g. posterior means and (co)variances
- Why are steps 2 and 3 hard?

Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior

data parameters

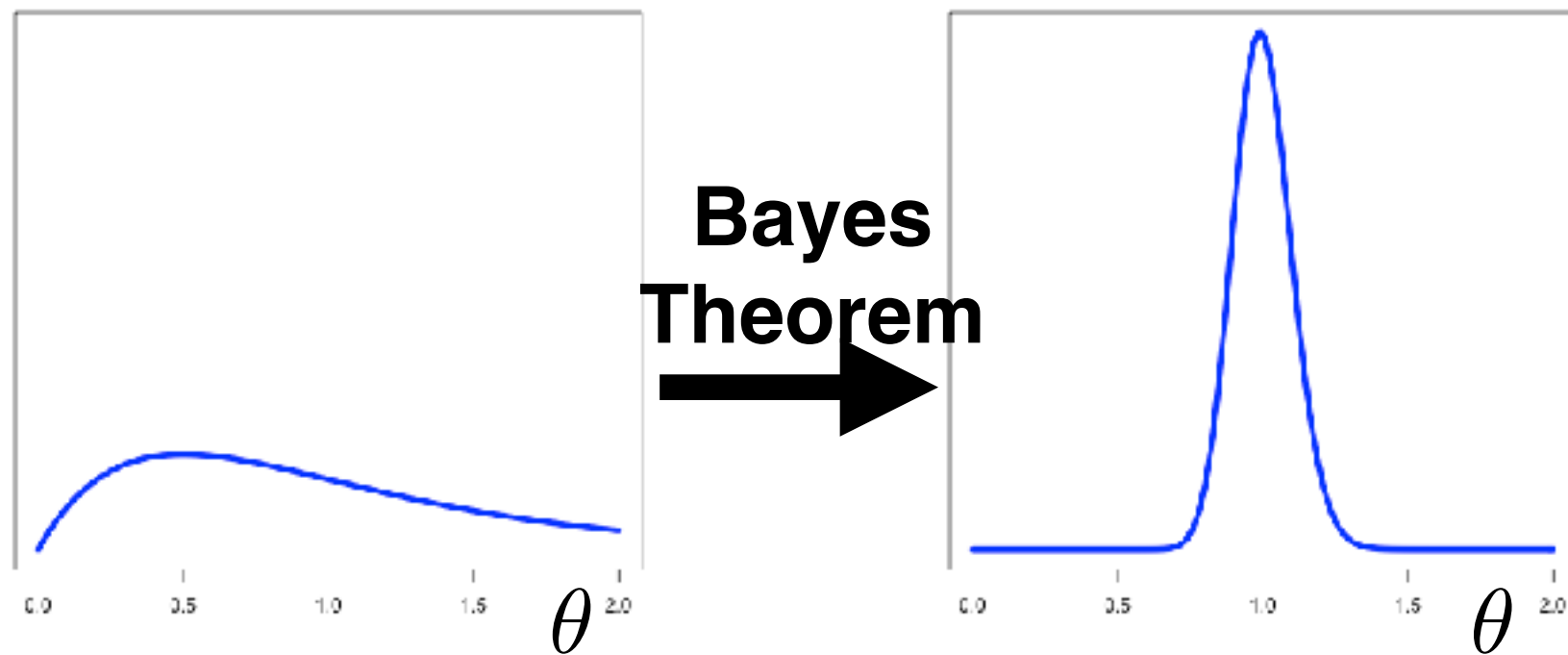


1. Build a model: choose prior & choose likelihood
 2. Compute the posterior
 3. Report a summary, e.g. posterior means and (co)variances
- Why are steps 2 and 3 hard?
 - Typically no closed form

Bayesian inference

$$p(\theta|y_{1:N}) \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

posterior likelihood prior

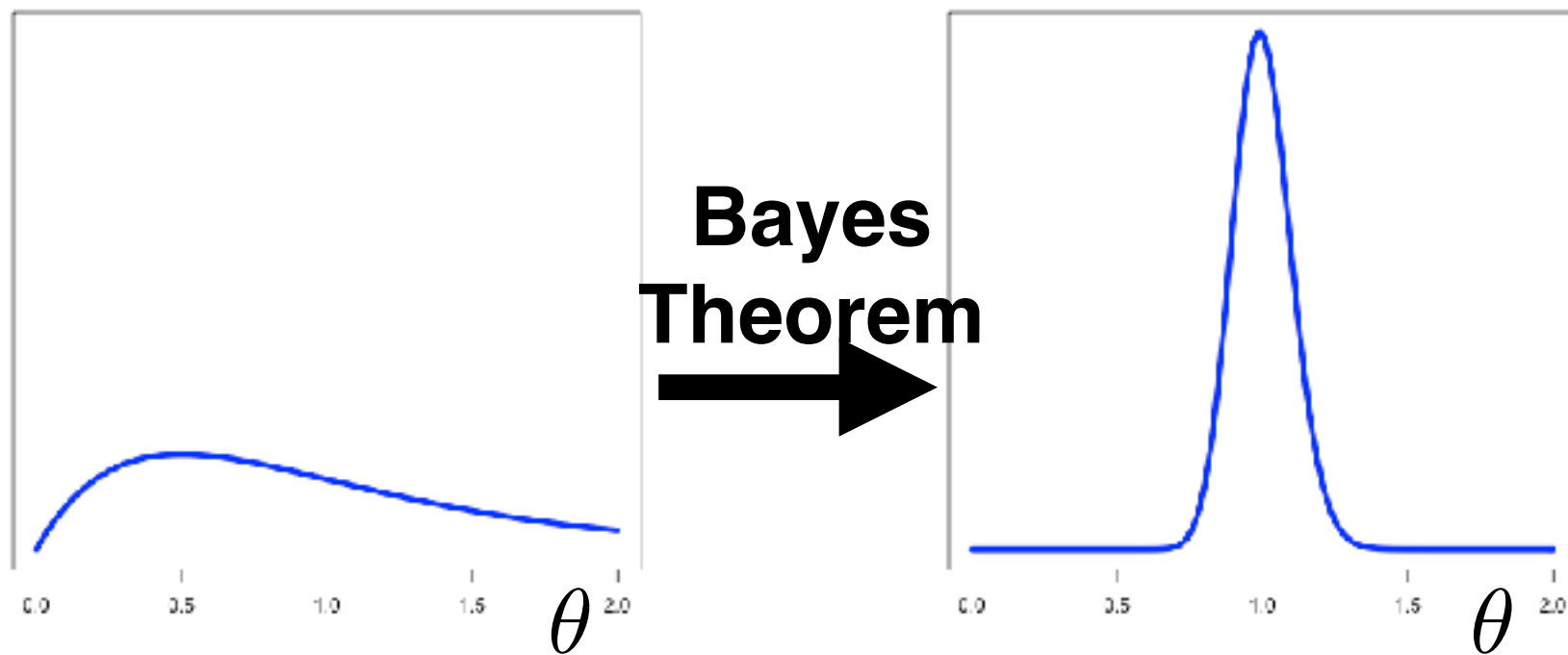


1. Build a model: choose prior & choose likelihood
 2. Compute the posterior
 3. Report a summary, e.g. posterior means and (co)variances
- Why are steps 2 and 3 hard?
 - Typically no closed form, high-dimensional integration

Bayesian inference

$$p(\theta|y_{1:N}) = p(y_{1:N}|\theta)p(\theta)/p(y_{1:N})$$

posterior likelihood prior

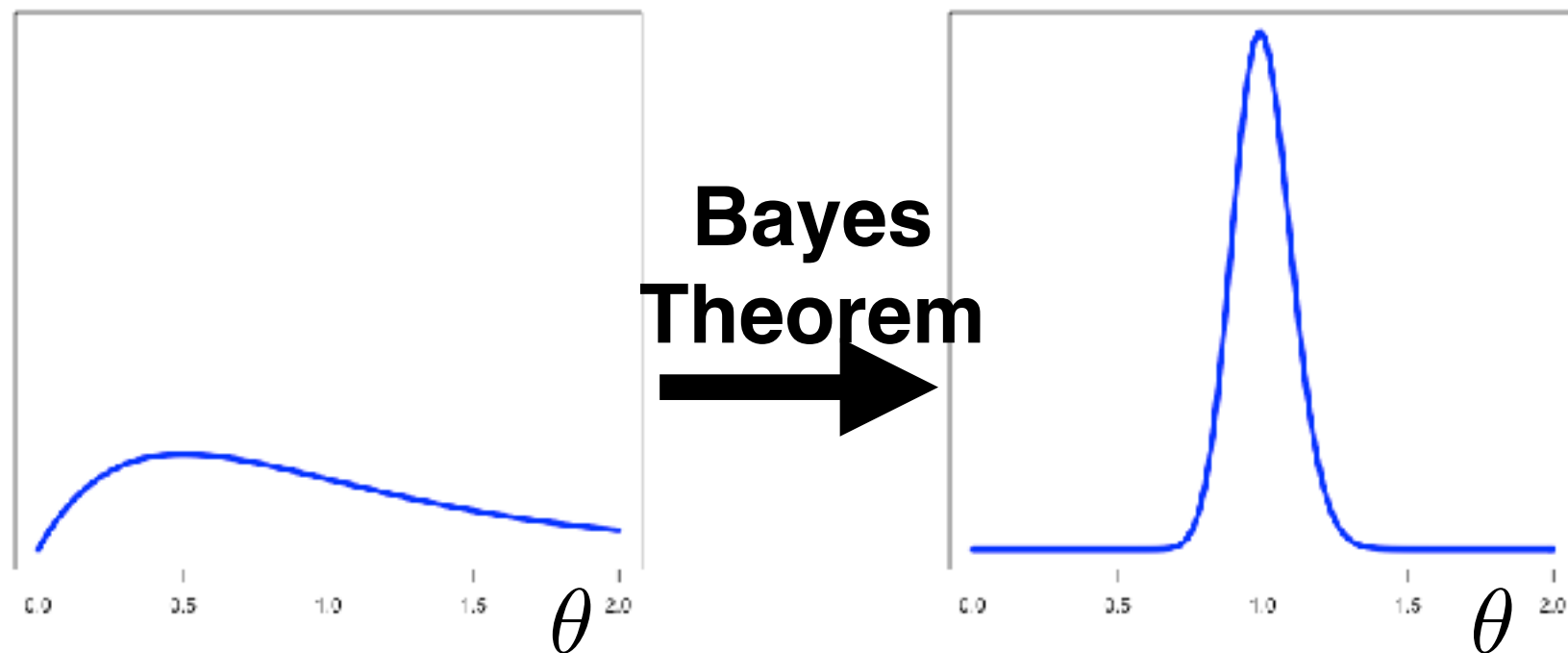


1. Build a model: choose prior & choose likelihood
 2. Compute the posterior
 3. Report a summary, e.g. posterior means and (co)variances
- Why are steps 2 and 3 hard?
 - Typically no closed form, high-dimensional integration

Bayesian inference

$$p(\theta|y_{1:N}) = p(y_{1:N}|\theta)p(\theta)/p(y_{1:N})$$

posterior likelihood prior evidence



1. Build a model: choose prior & choose likelihood
 2. Compute the posterior
 3. Report a summary, e.g. posterior means and (co)variances
- Why are steps 2 and 3 hard?
 - Typically no closed form, high-dimensional integration

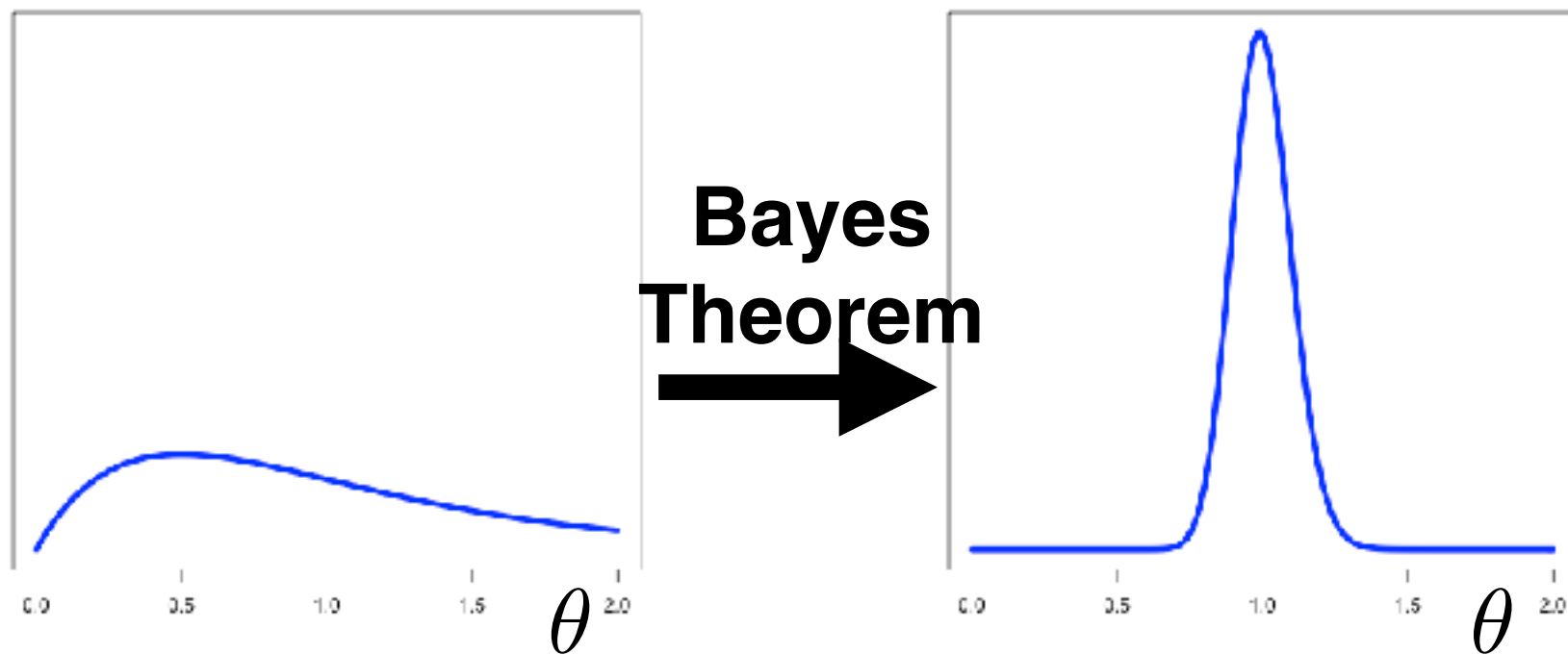
Bayesian inference

data parameters

↓ ↓

$$p(\theta|y_{1:N}) = p(y_{1:N}|\theta)p(\theta) / \int p(y_{1:N}, \theta) d\theta$$

posterior likelihood prior evidence



1. Build a model: choose prior & choose likelihood
 2. Compute the posterior
 3. Report a summary, e.g. posterior means and (co)variances
- Why are steps 2 and 3 hard?
 - Typically no closed form, high-dimensional integration

Approximate Bayesian Inference

Approximate Bayesian Inference

- Gold standard: Markov Chain Monte Carlo (MCMC)

[Bardenet,
Doucet,
Holmes
2017]

Approximate Bayesian Inference

- Gold standard: Markov Chain Monte Carlo (MCMC)
 - Eventually accurate but can be slow

[Bardenet,
Doucet,
Holmes
2017]

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow

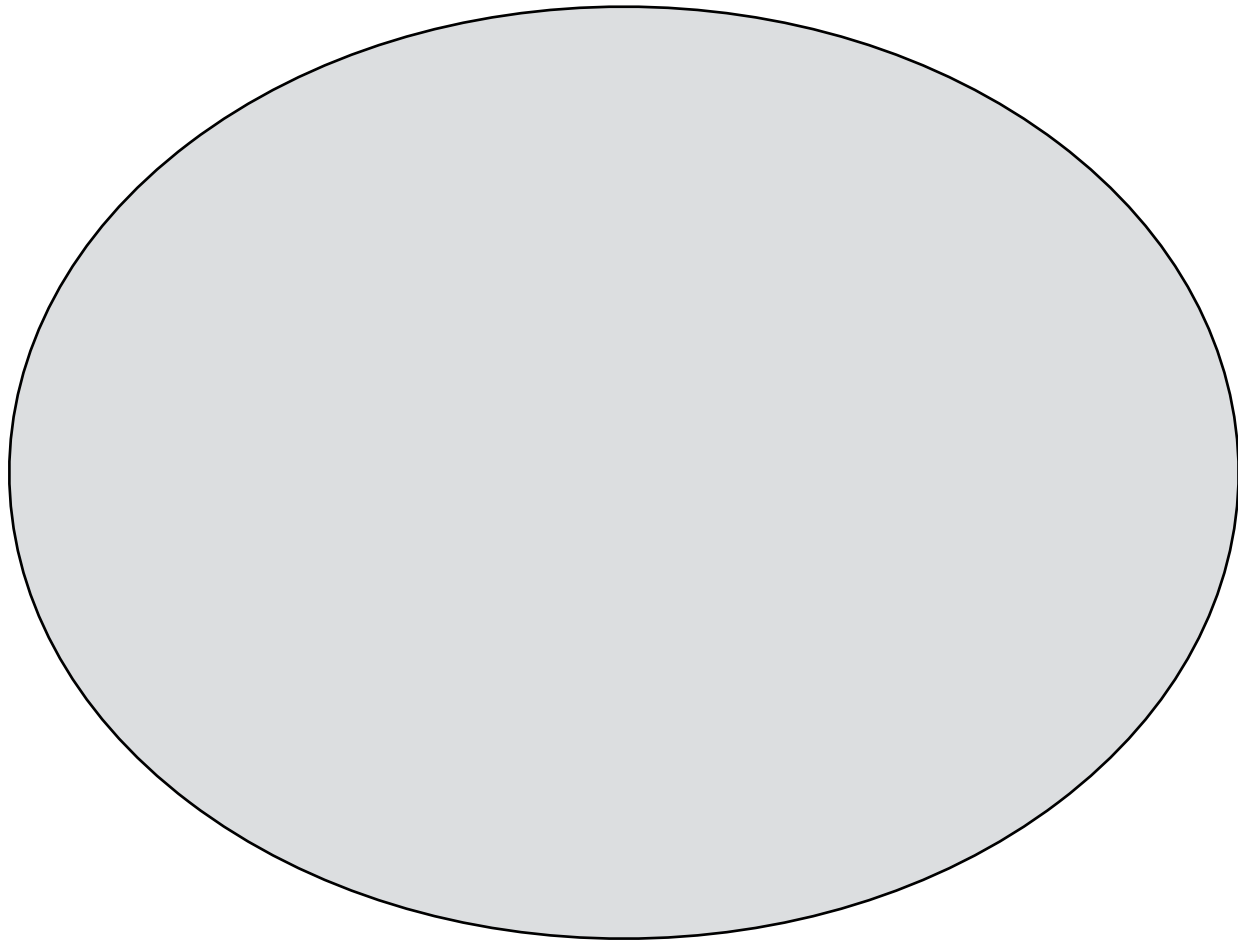
Instead: an optimization approach

- Approximate posterior with q^*

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
 - Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

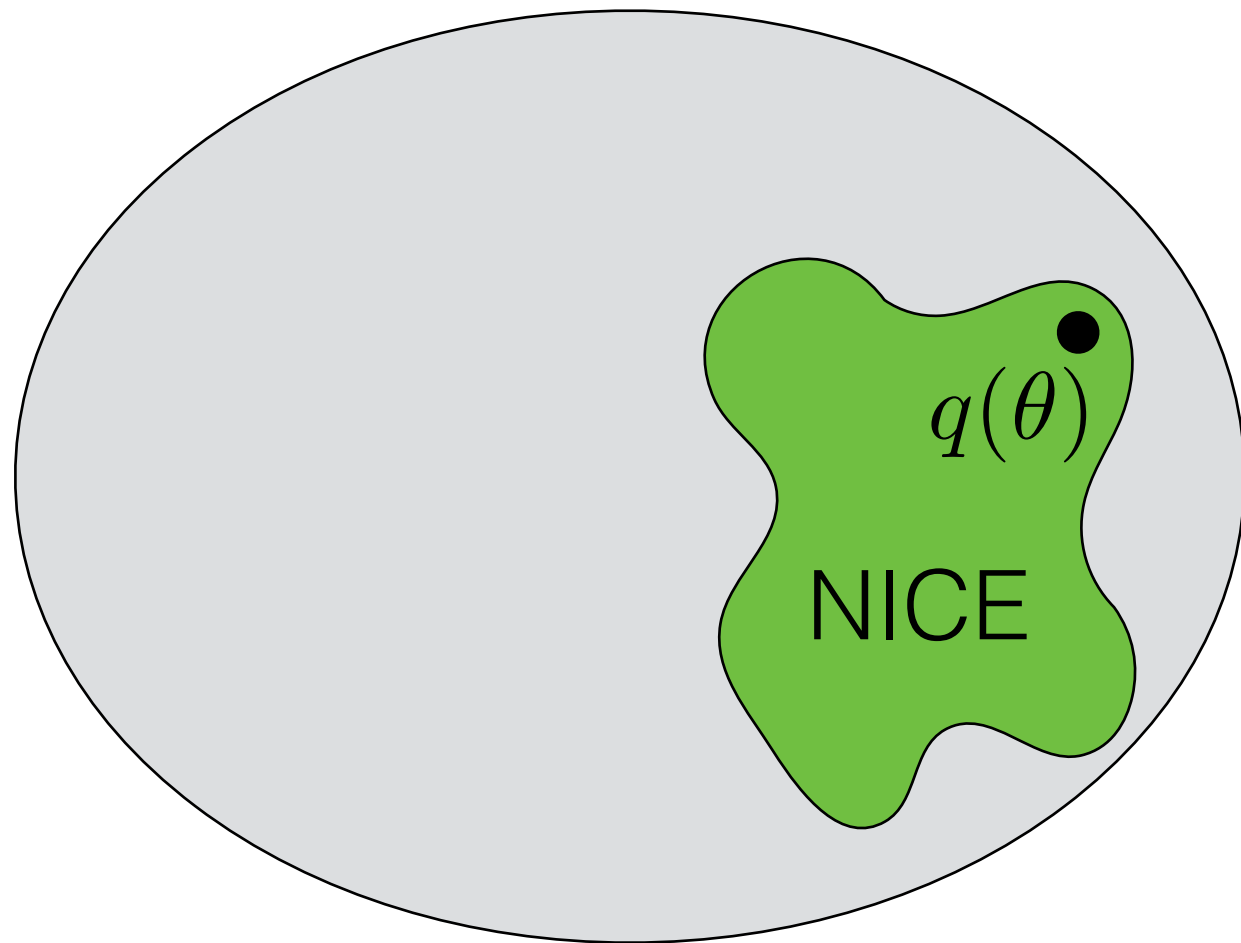
Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
 - Eventually accurate but can be slow

Instead: an optimization approach

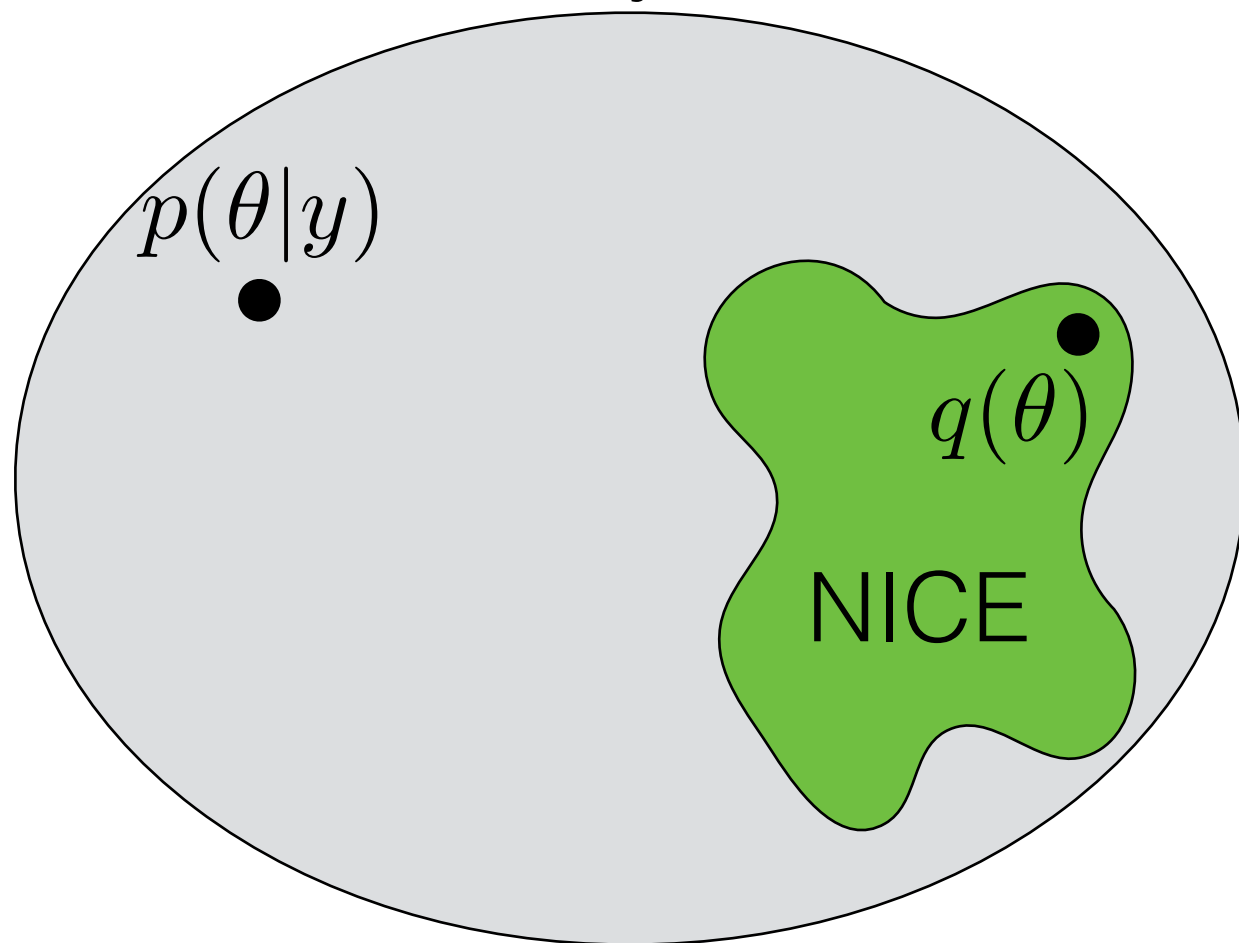
- Approximate posterior with q^*



Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
 - Eventually accurate but can be slow



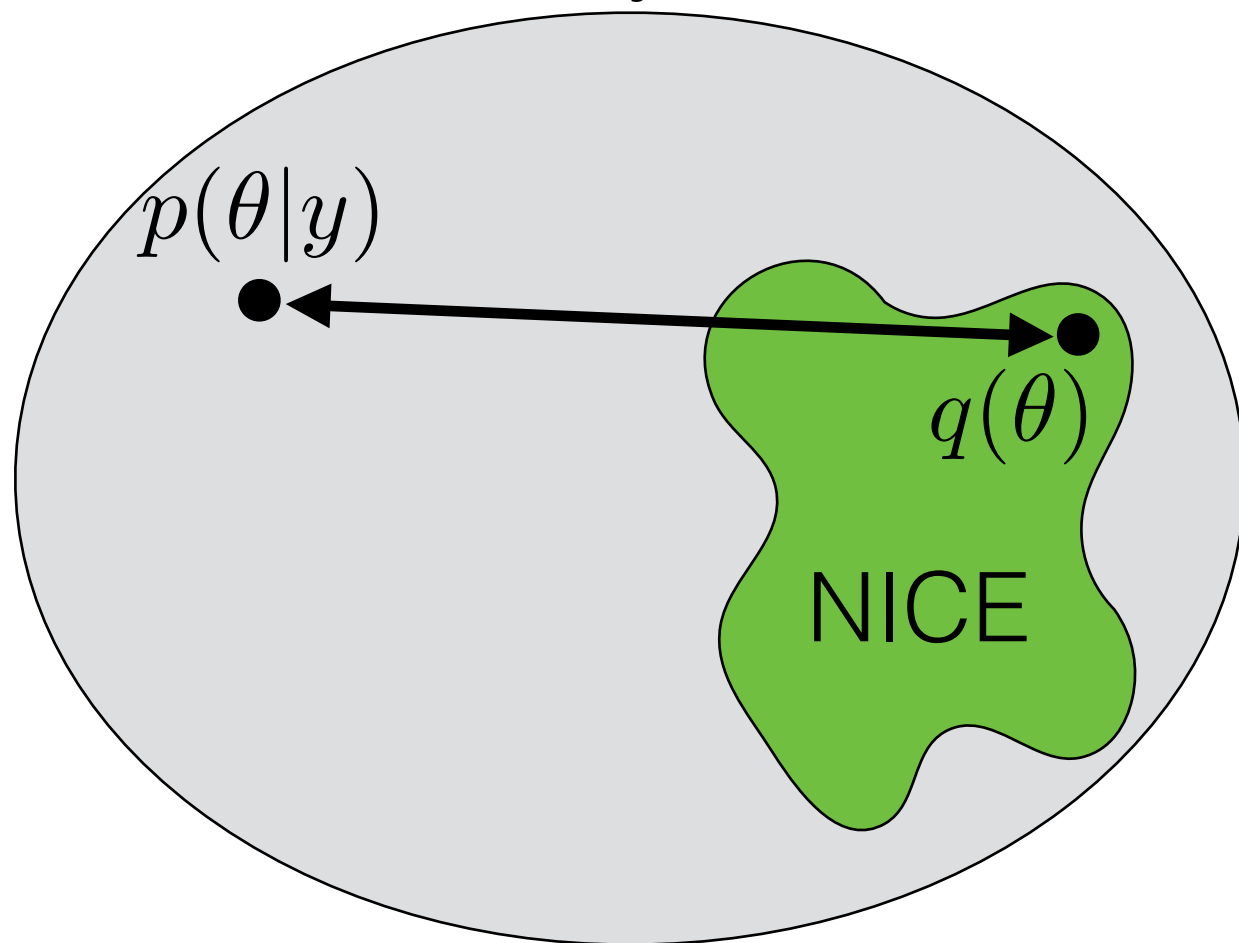
Instead: an optimization approach

- Approximate posterior with q^*

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
 - Eventually accurate but can be slow



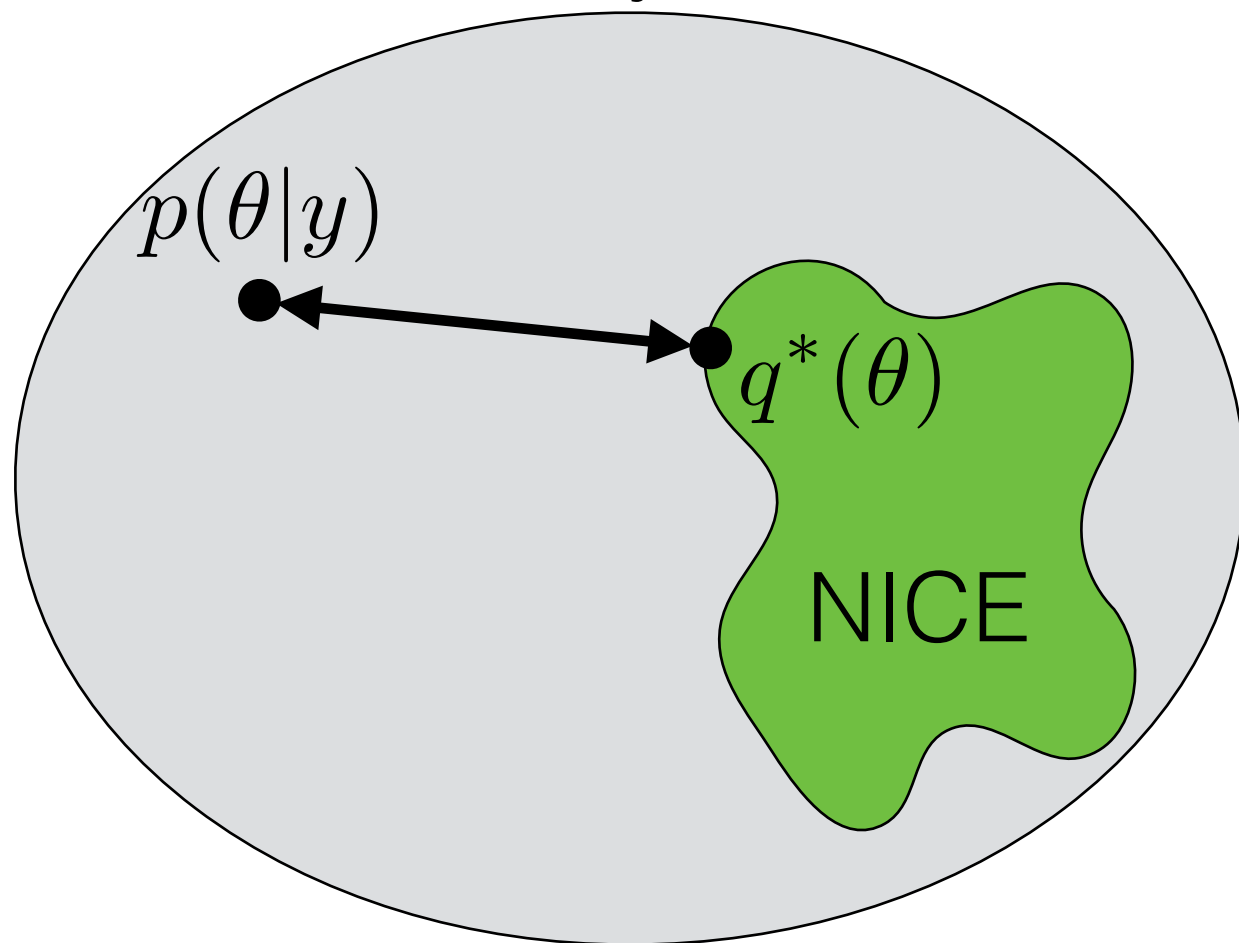
Instead: an optimization approach

- Approximate posterior with q^*

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
 - Eventually accurate but can be slow



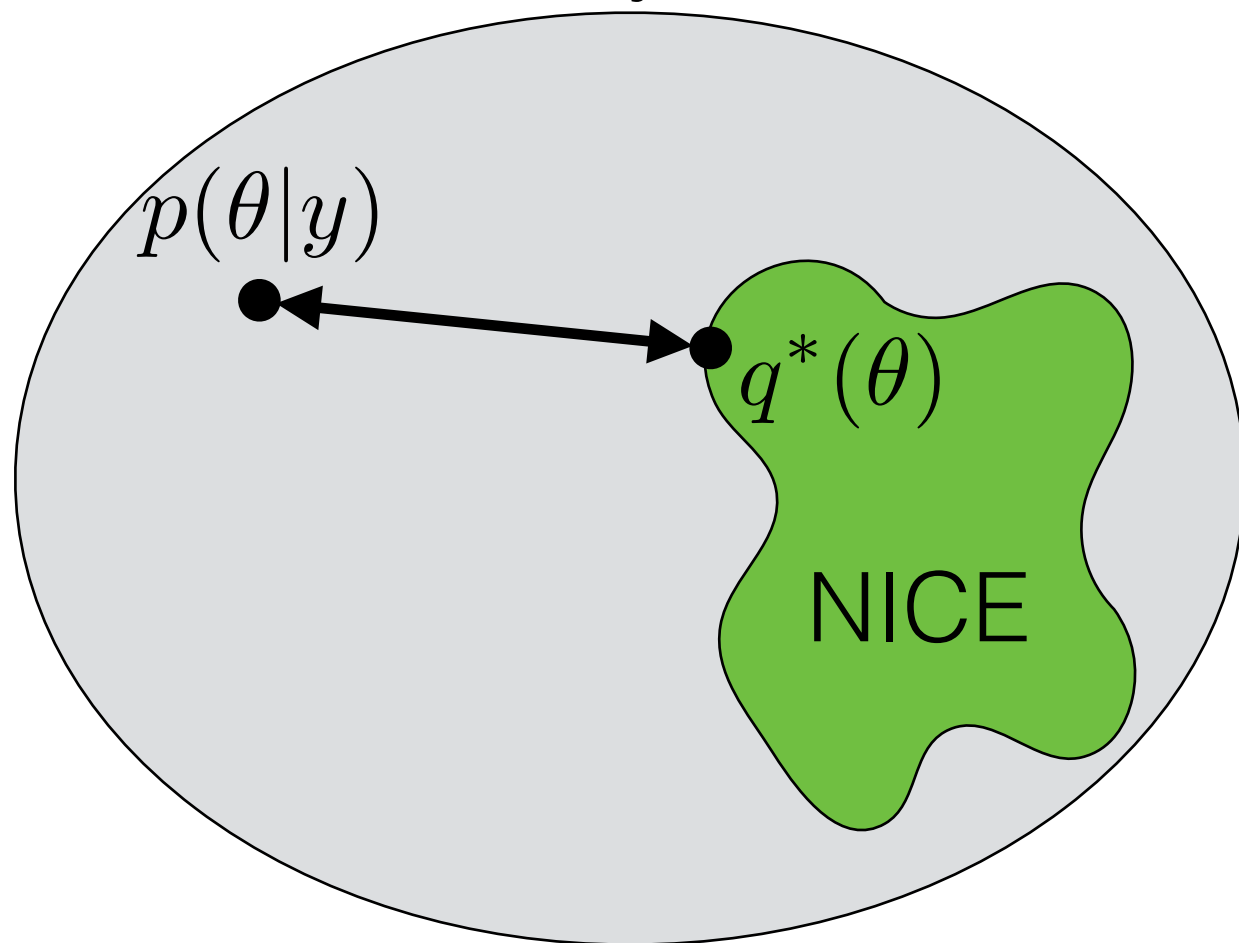
Instead: an optimization approach

- Approximate posterior with q^*

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

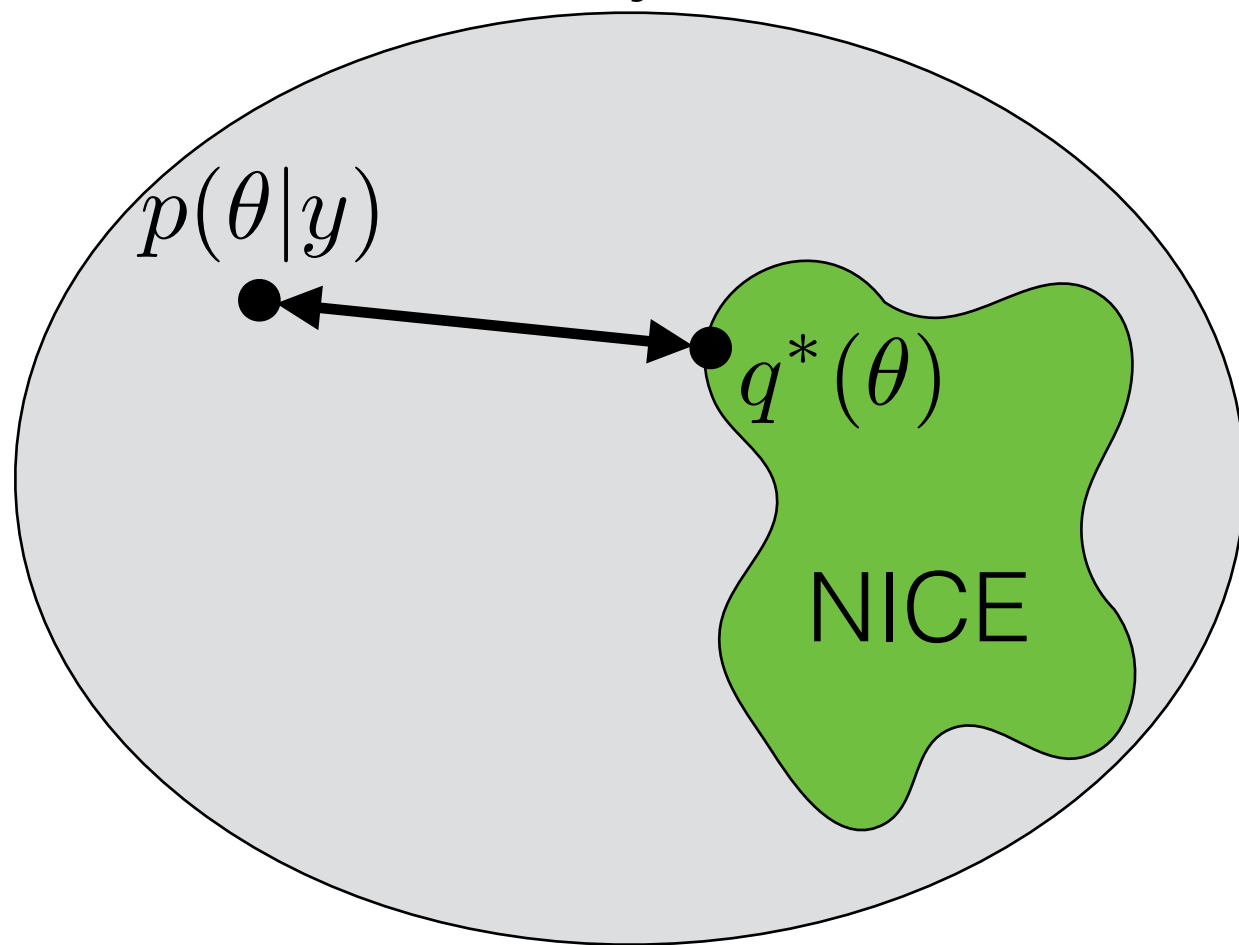
- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

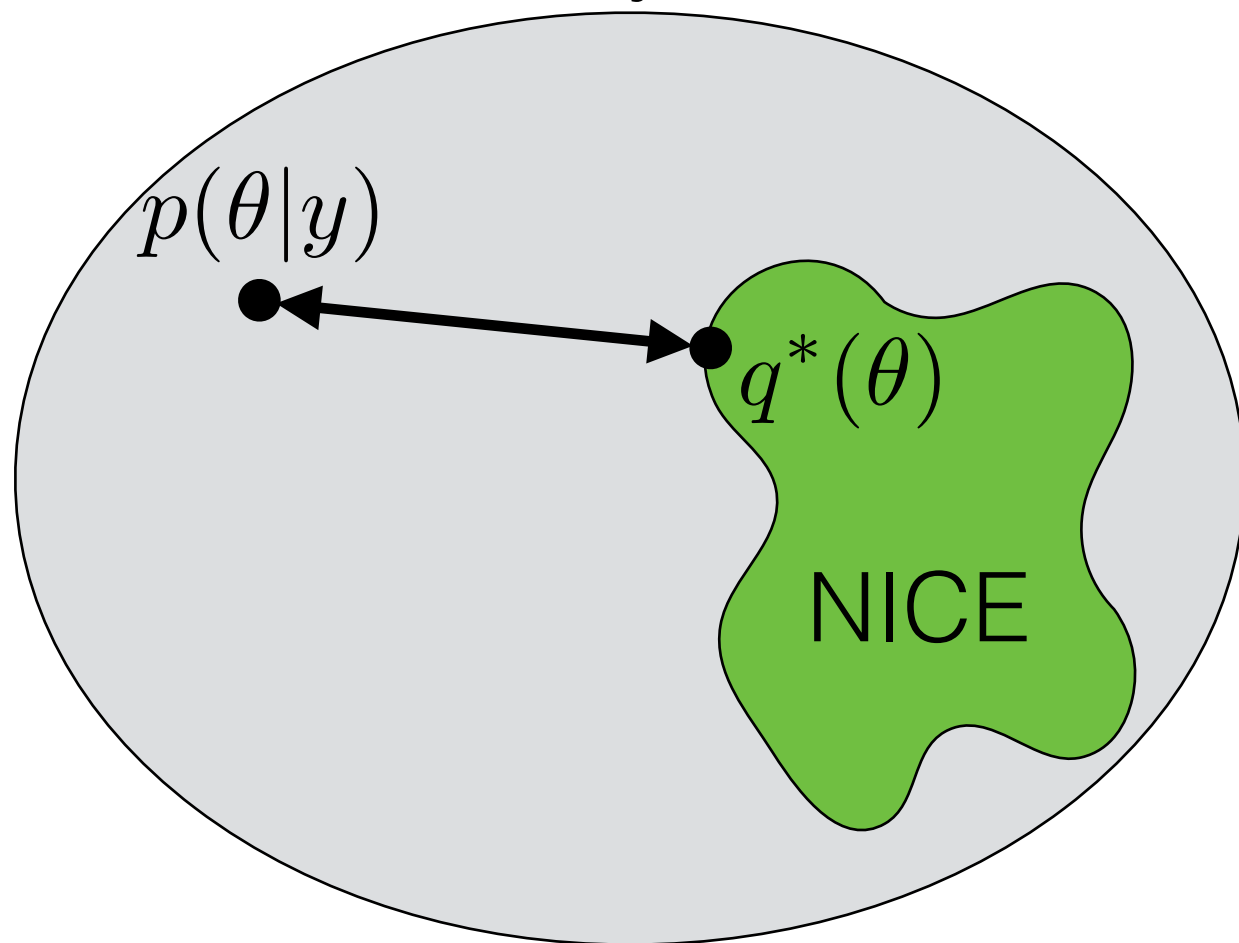
- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

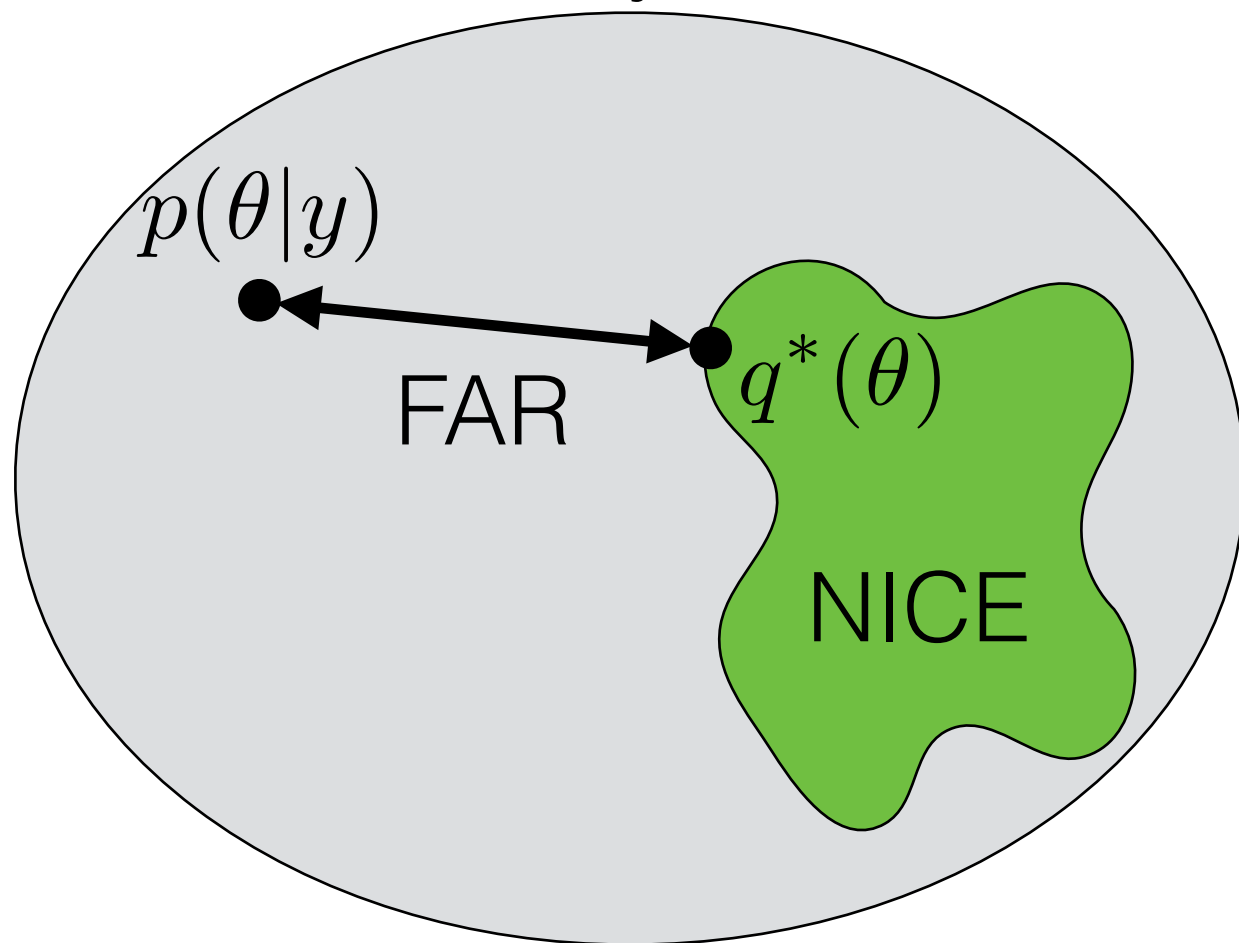
- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

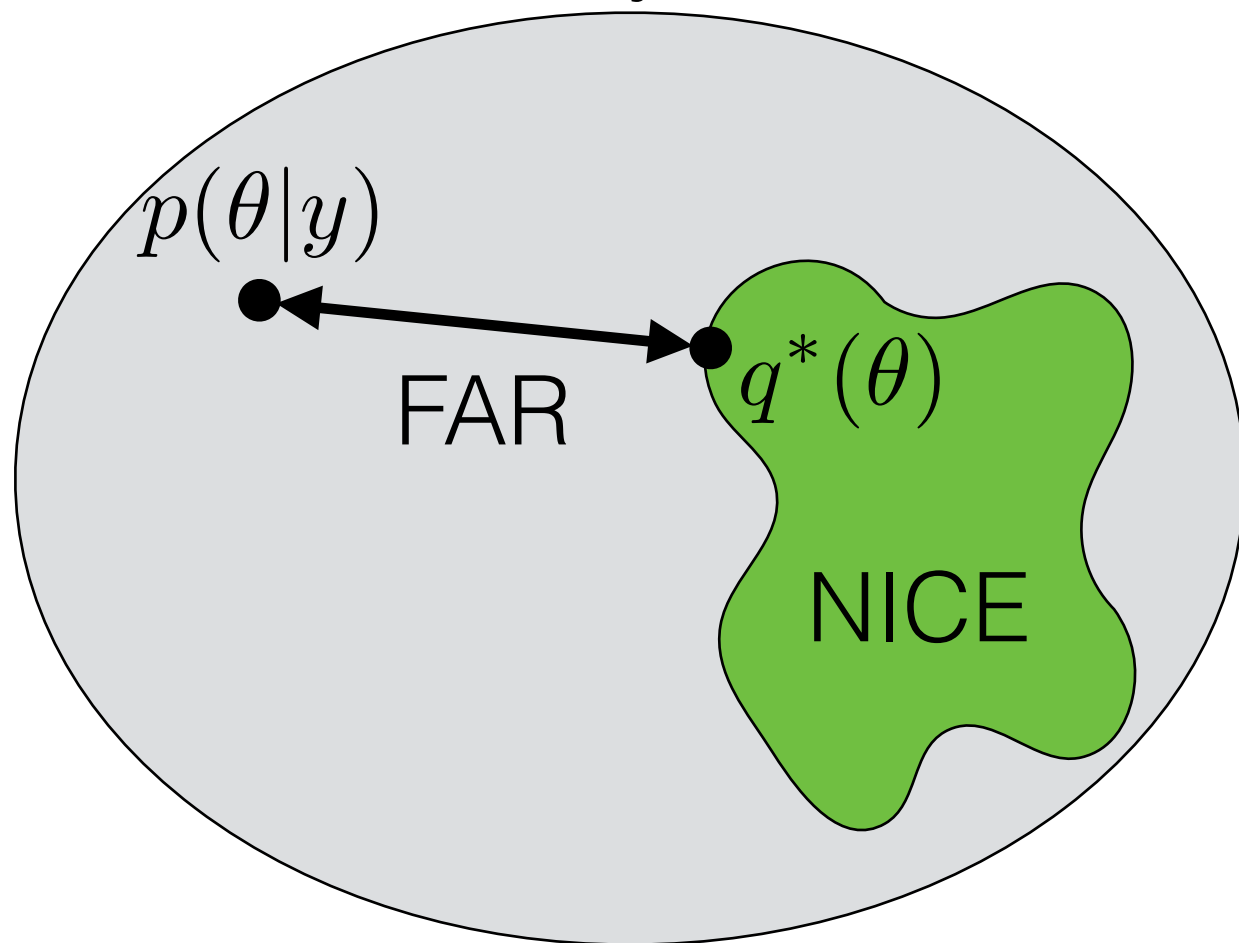
- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

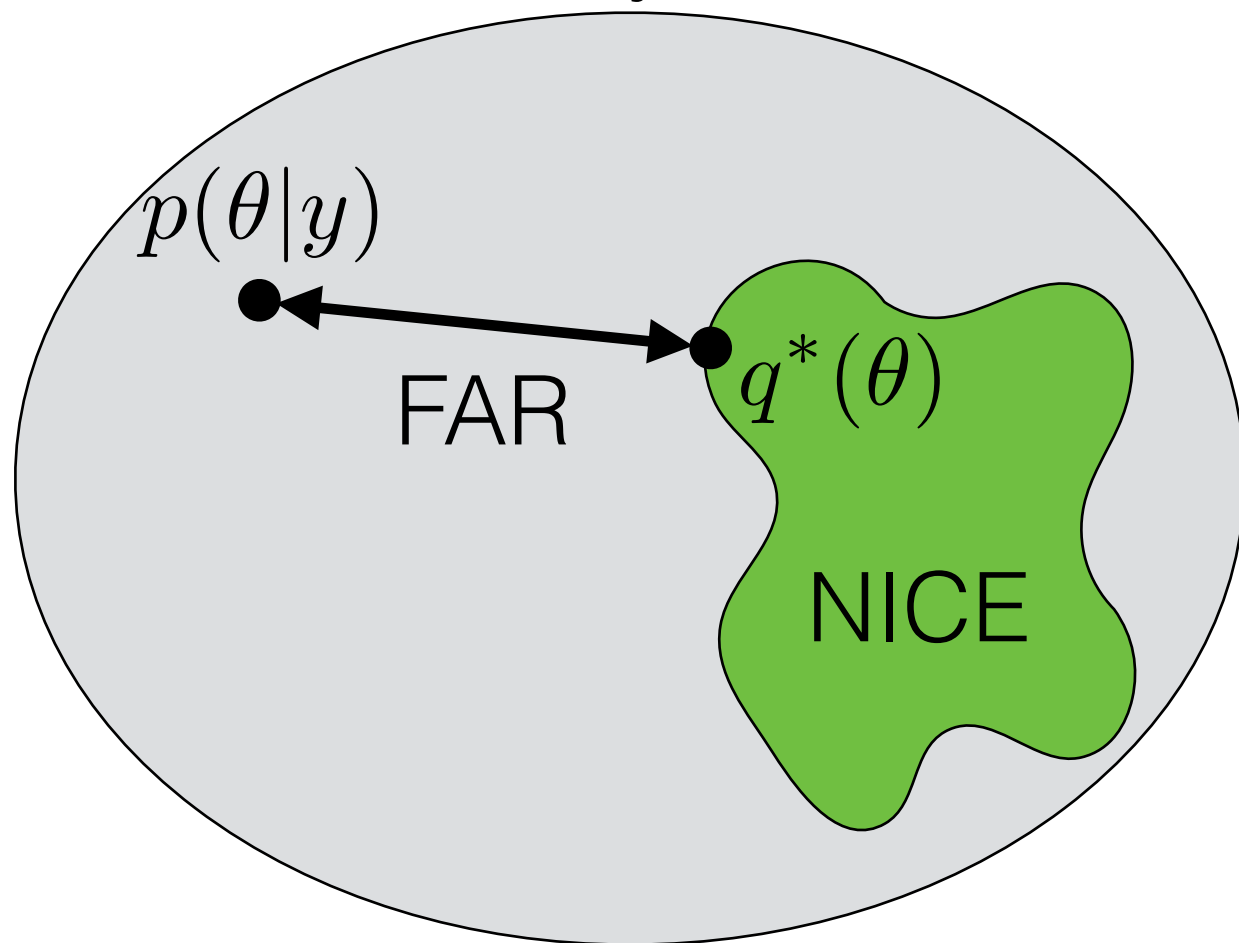
- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

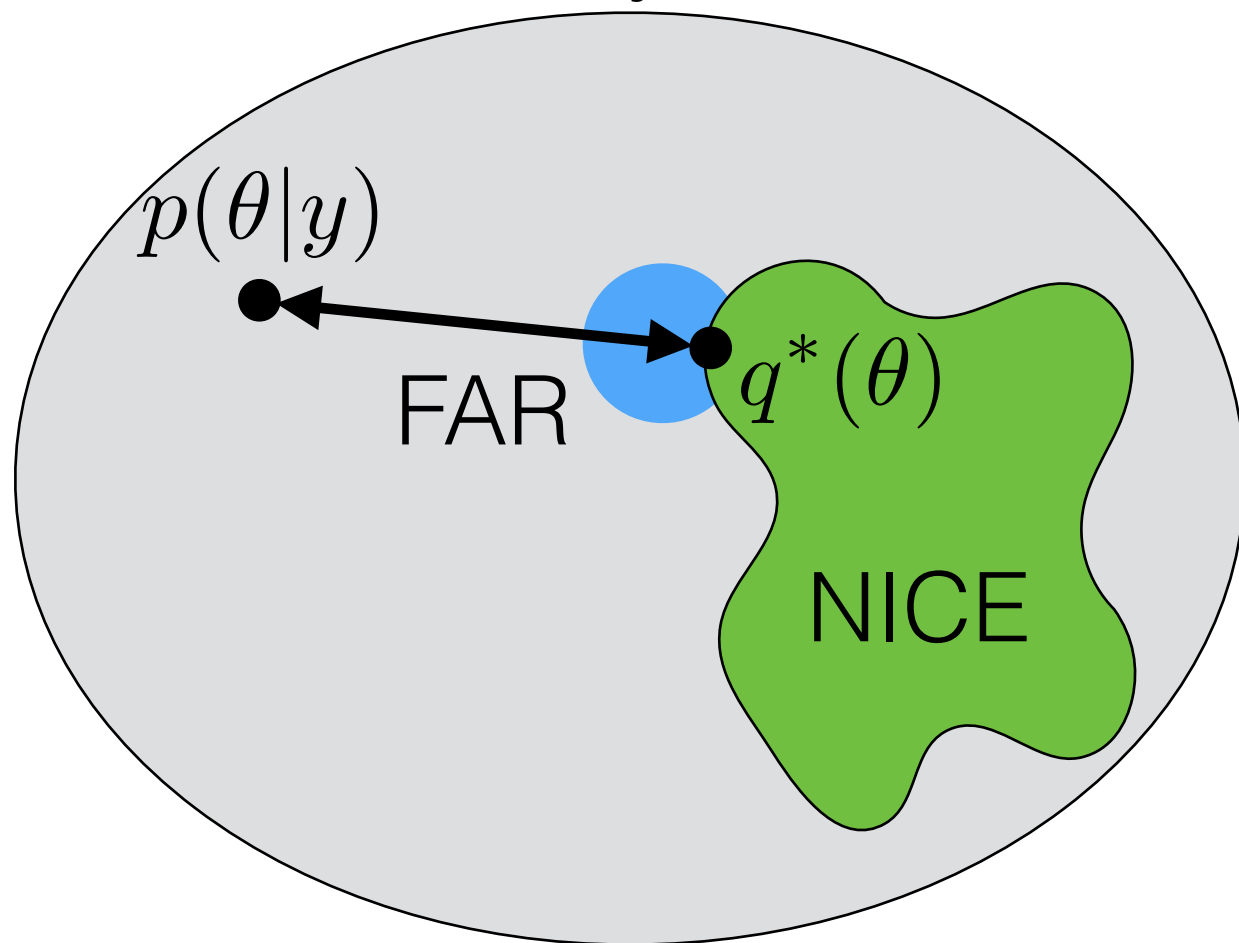
- Variational Bayes (VB): f is Kullback-Leibler divergence

$$KL(q(\cdot) || p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

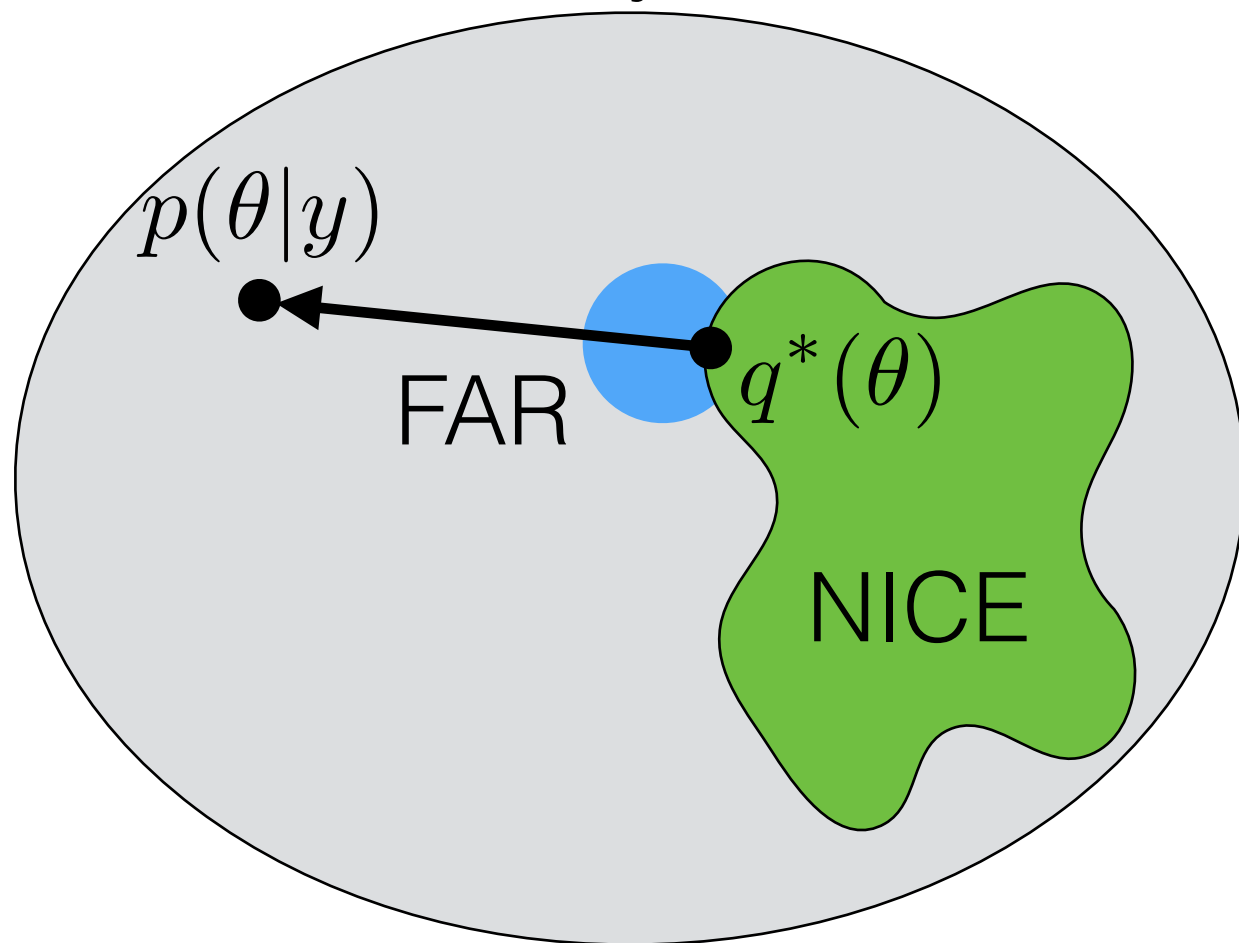
- Variational Bayes (VB): f is Kullback-Leibler divergence

$$KL(q(\cdot) || p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

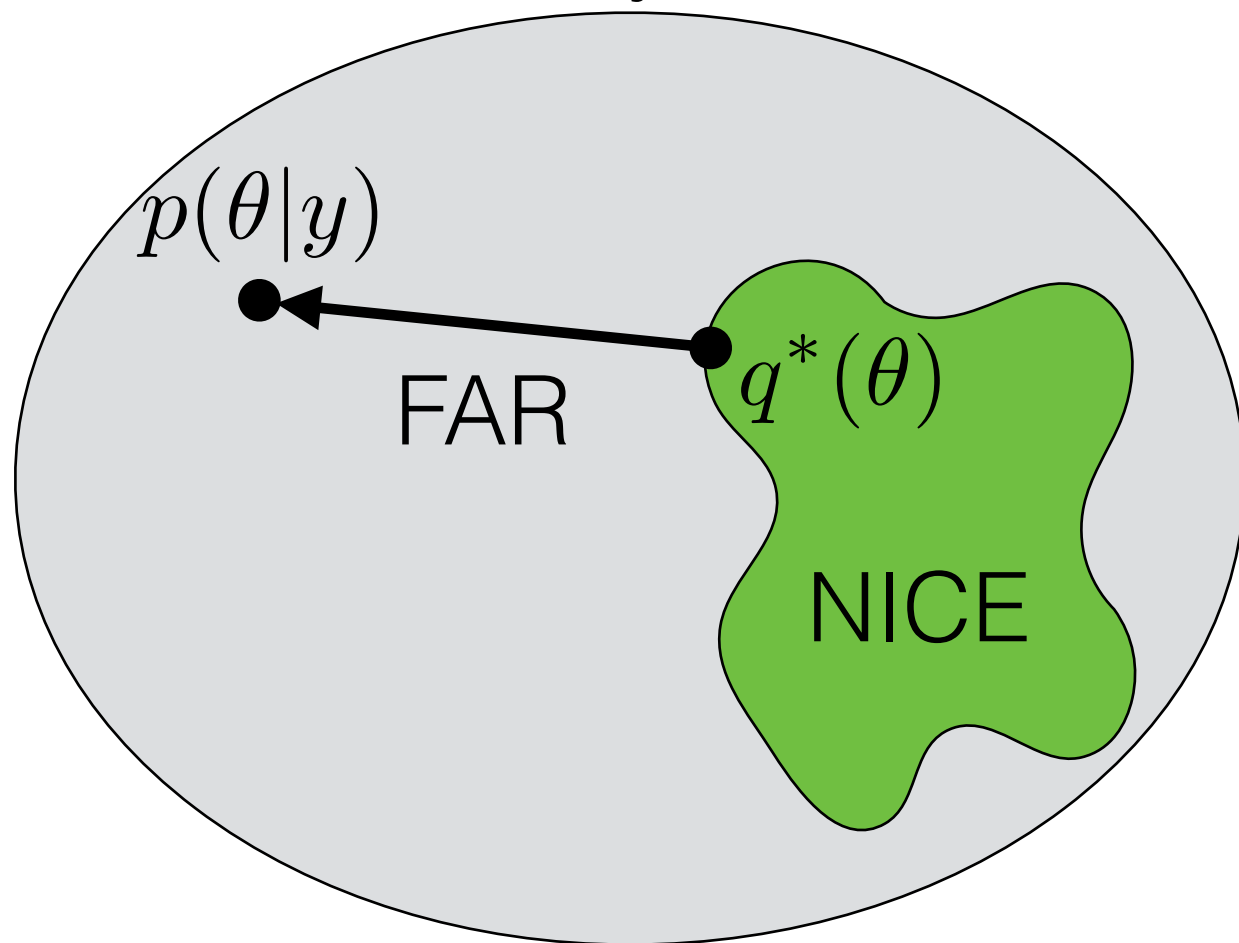
- Variational Bayes (VB): f is Kullback-Leibler divergence

$$KL(q(\cdot) || p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

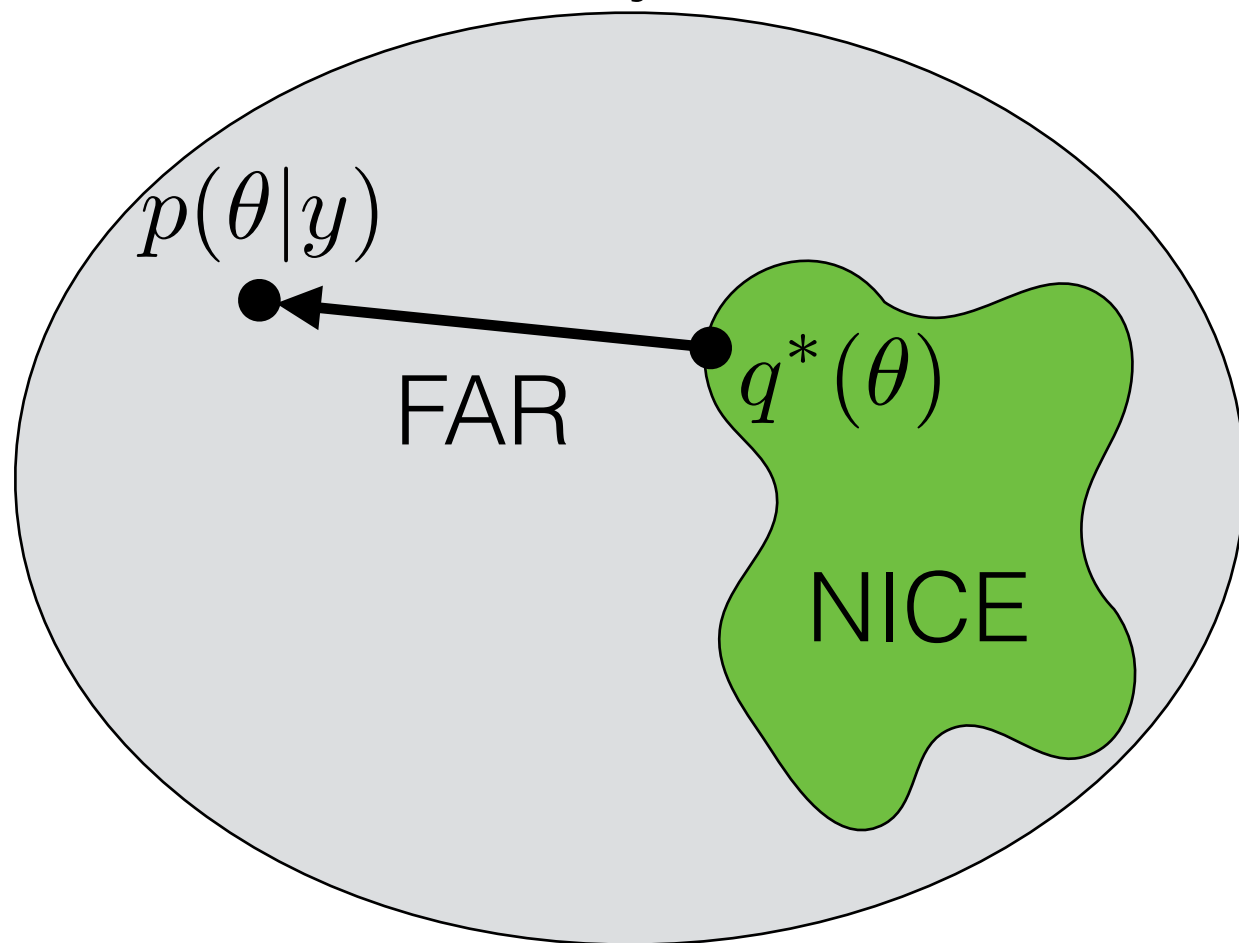
- Variational Bayes (VB): f is Kullback-Leibler divergence

$$KL(q(\cdot) || p(\cdot|y))$$

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
 - Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

- Variational Bayes (VB): f is Kullback-Leibler divergence

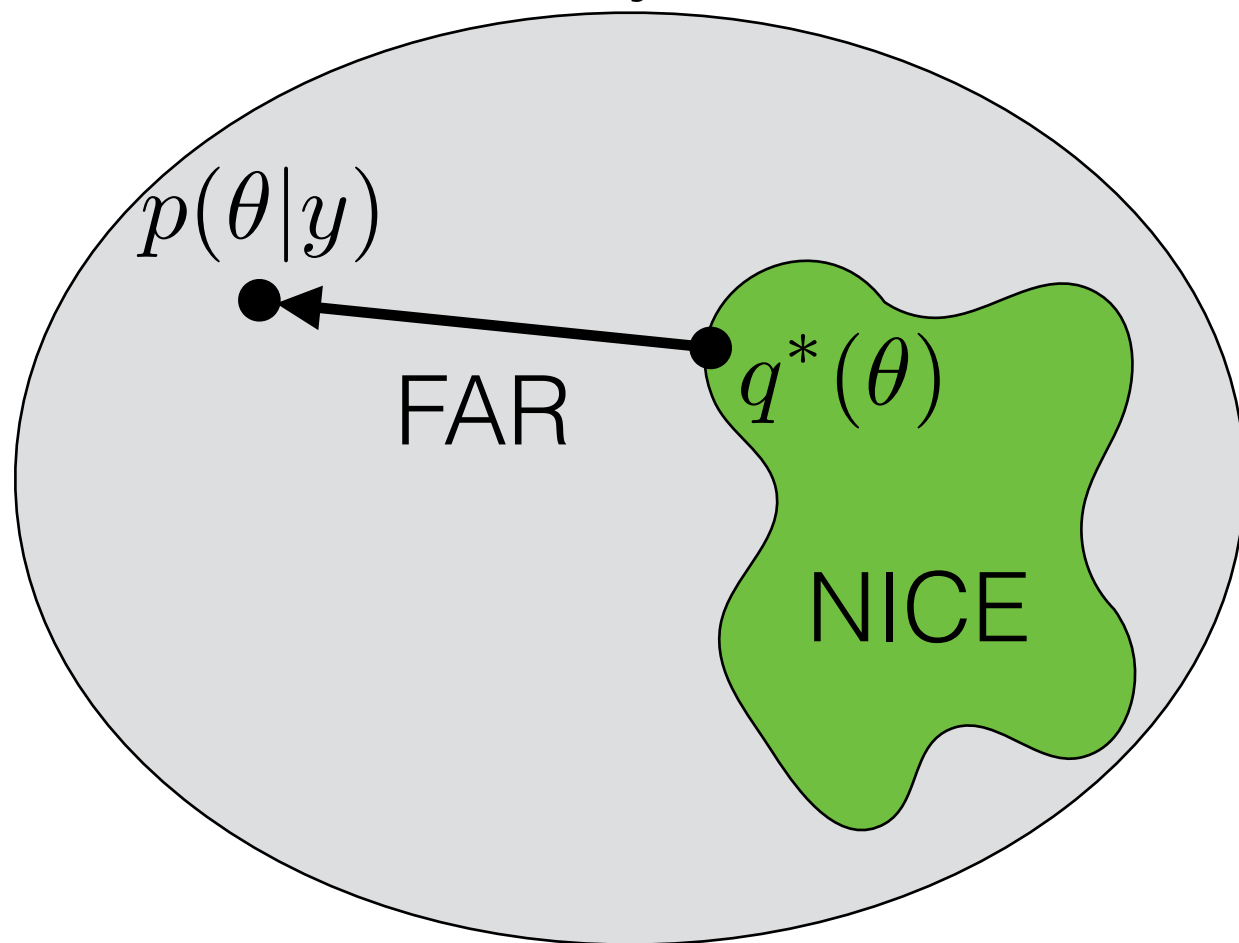
$$KL(q(\cdot) || p(\cdot|y))$$

- VB practical success

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

- Variational Bayes (VB): f is Kullback-Leibler divergence

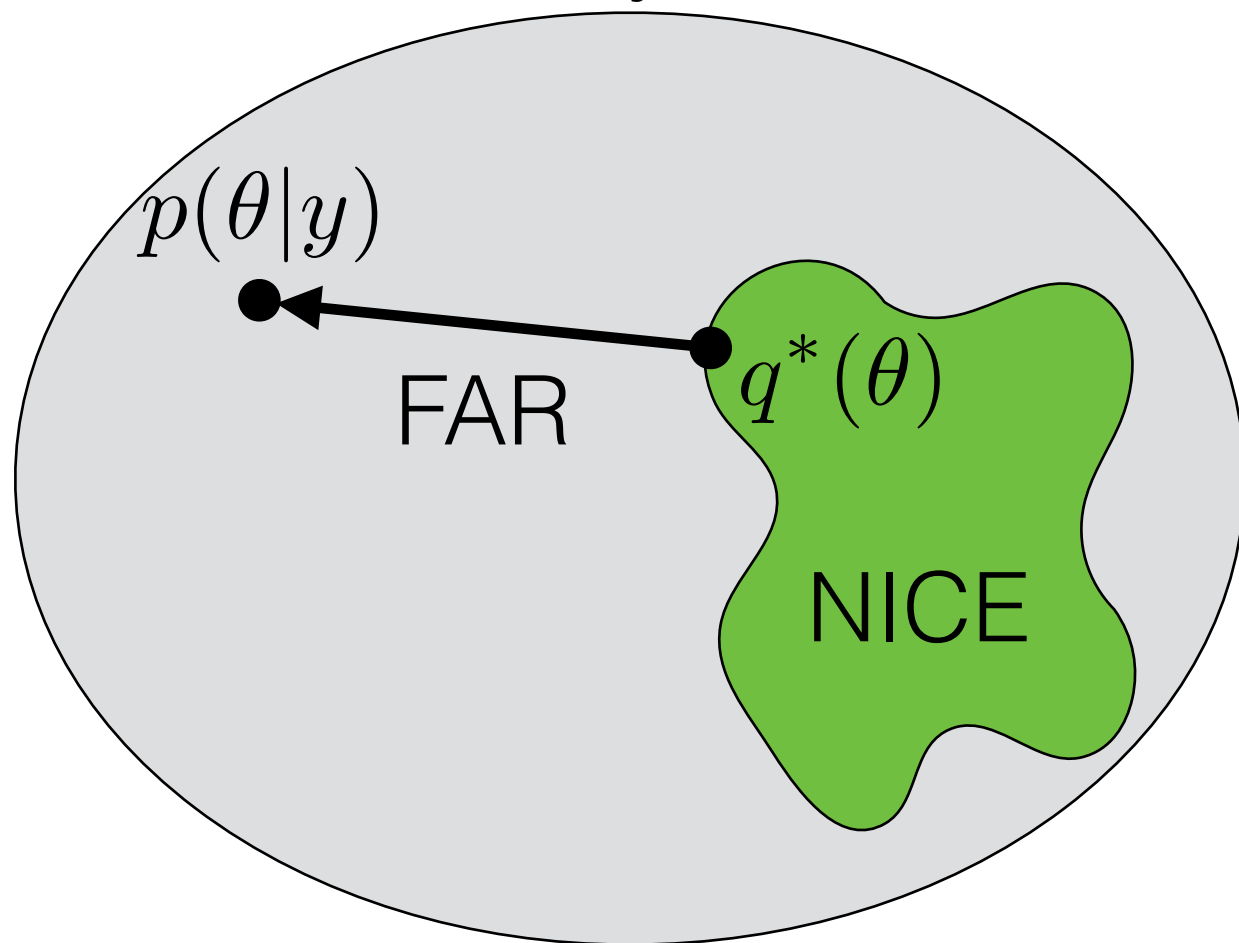
$$KL(q(\cdot) || p(\cdot|y))$$

- VB practical success: point estimates and prediction

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

- Variational Bayes (VB): f is Kullback-Leibler divergence

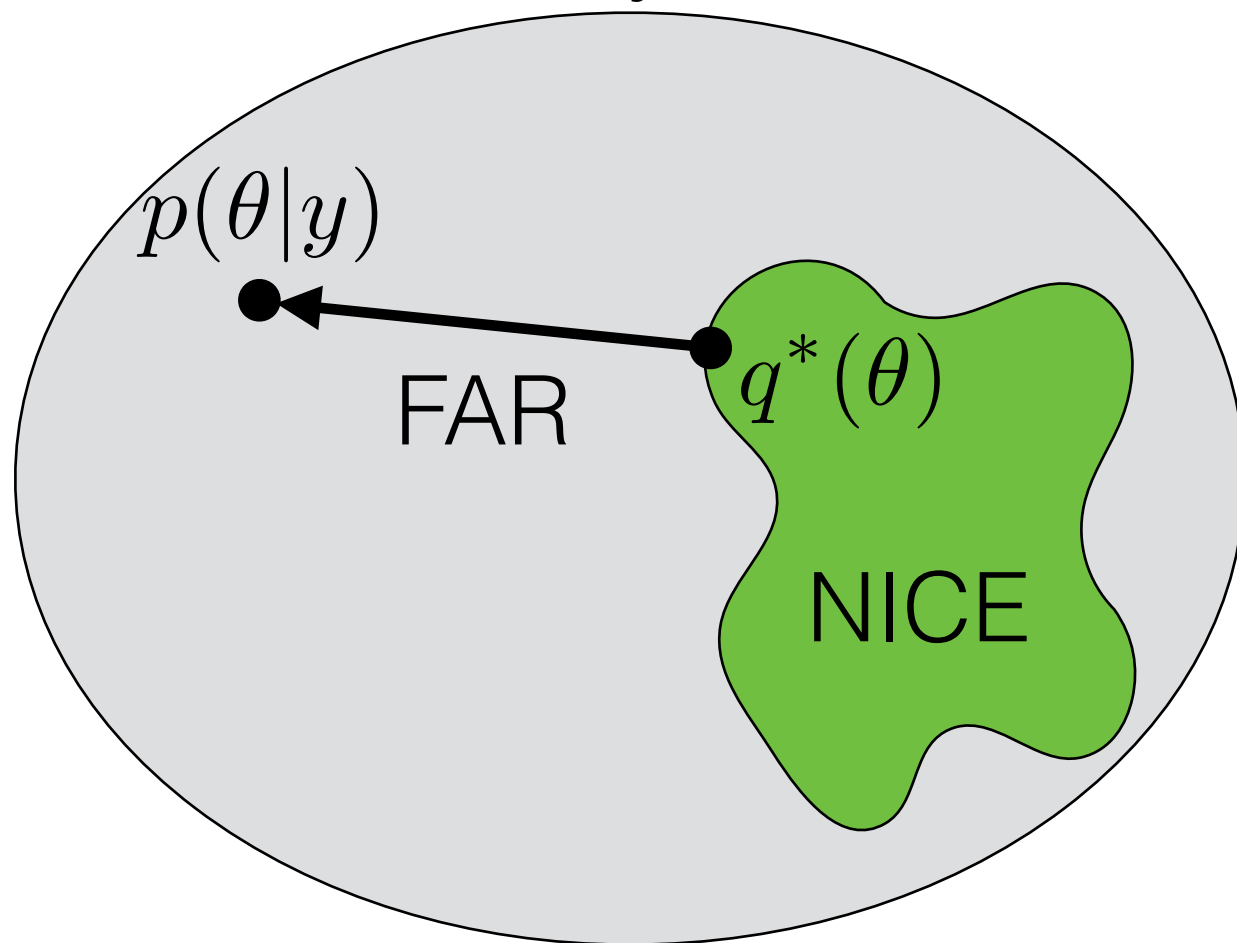
$$KL(q(\cdot) || p(\cdot|y))$$

- VB practical success: point estimates and prediction, fast

Approximate Bayesian Inference

[Bardenet,
Doucet,
Holmes
2017]

- Gold standard: Markov Chain Monte Carlo (MCMC)
- Eventually accurate but can be slow



Instead: an optimization approach

- Approximate posterior with q^*

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

- Variational Bayes (VB): f is Kullback-Leibler divergence

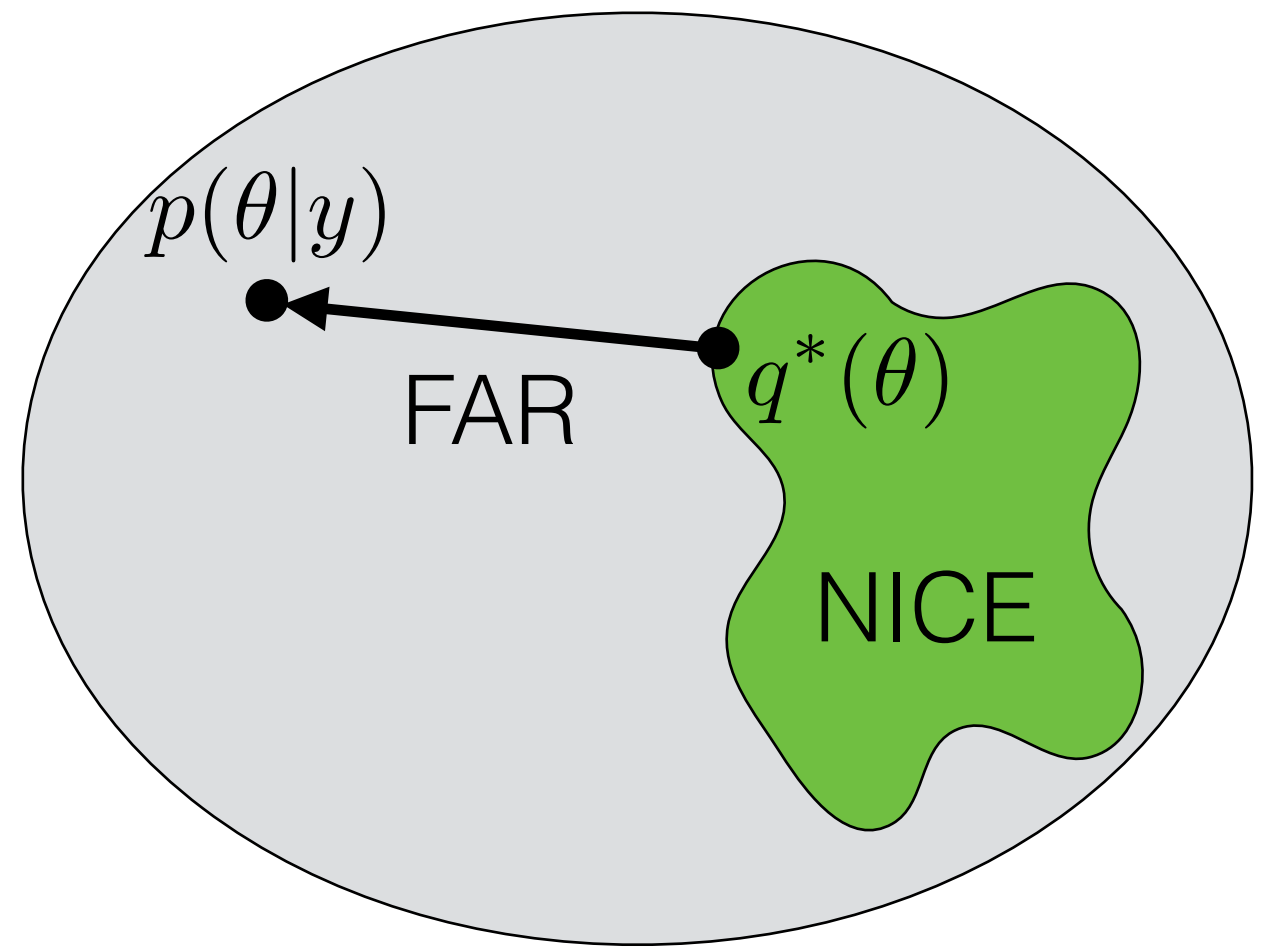
$$KL(q(\cdot) || p(\cdot|y))$$

- VB practical success: point estimates and prediction, fast, streaming, distributed (3.6M Wikipedia, 350K Nature)

Why KL?

- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$



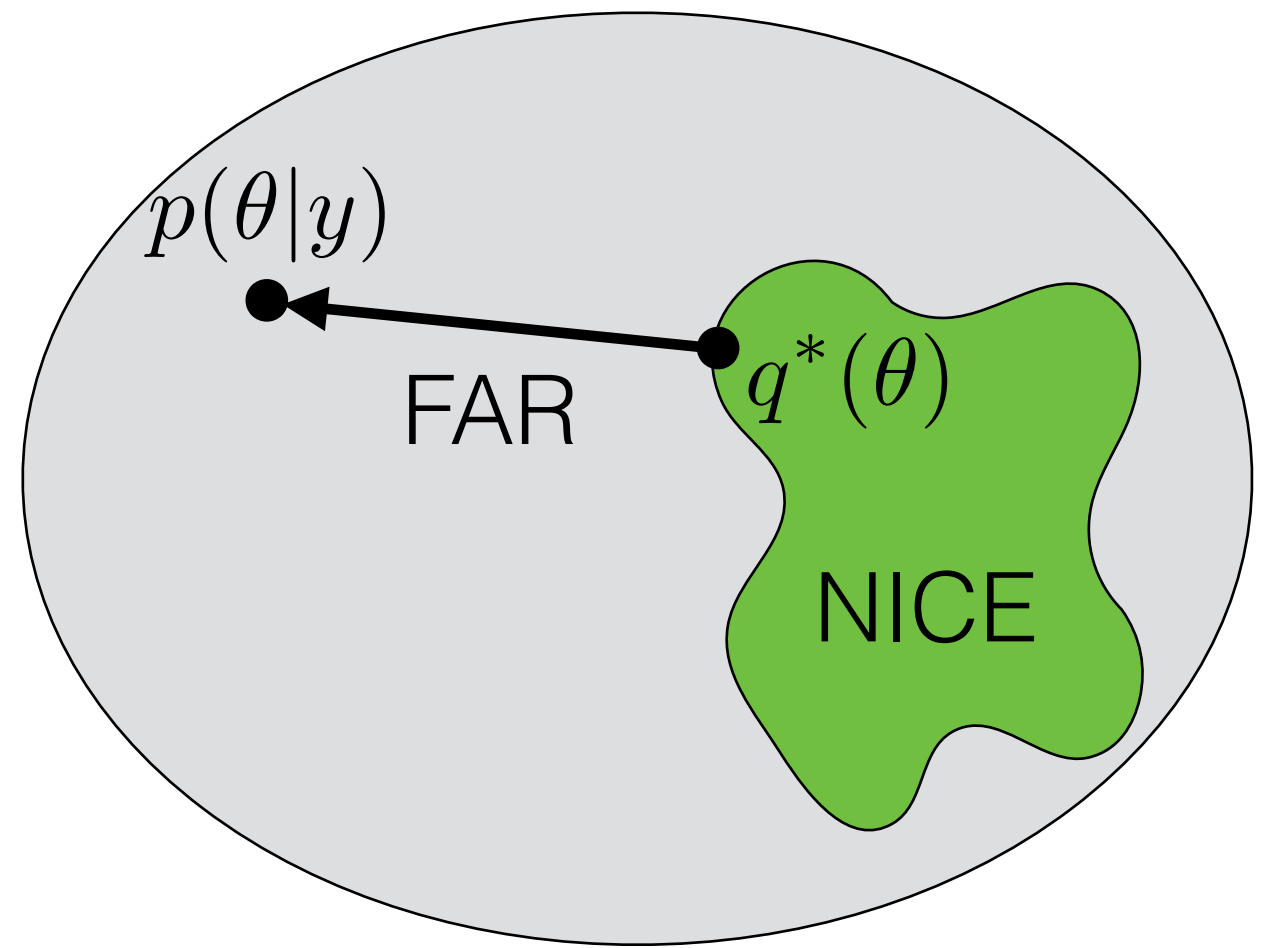
Why KL?

- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL} (q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL} (q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$



Why KL?

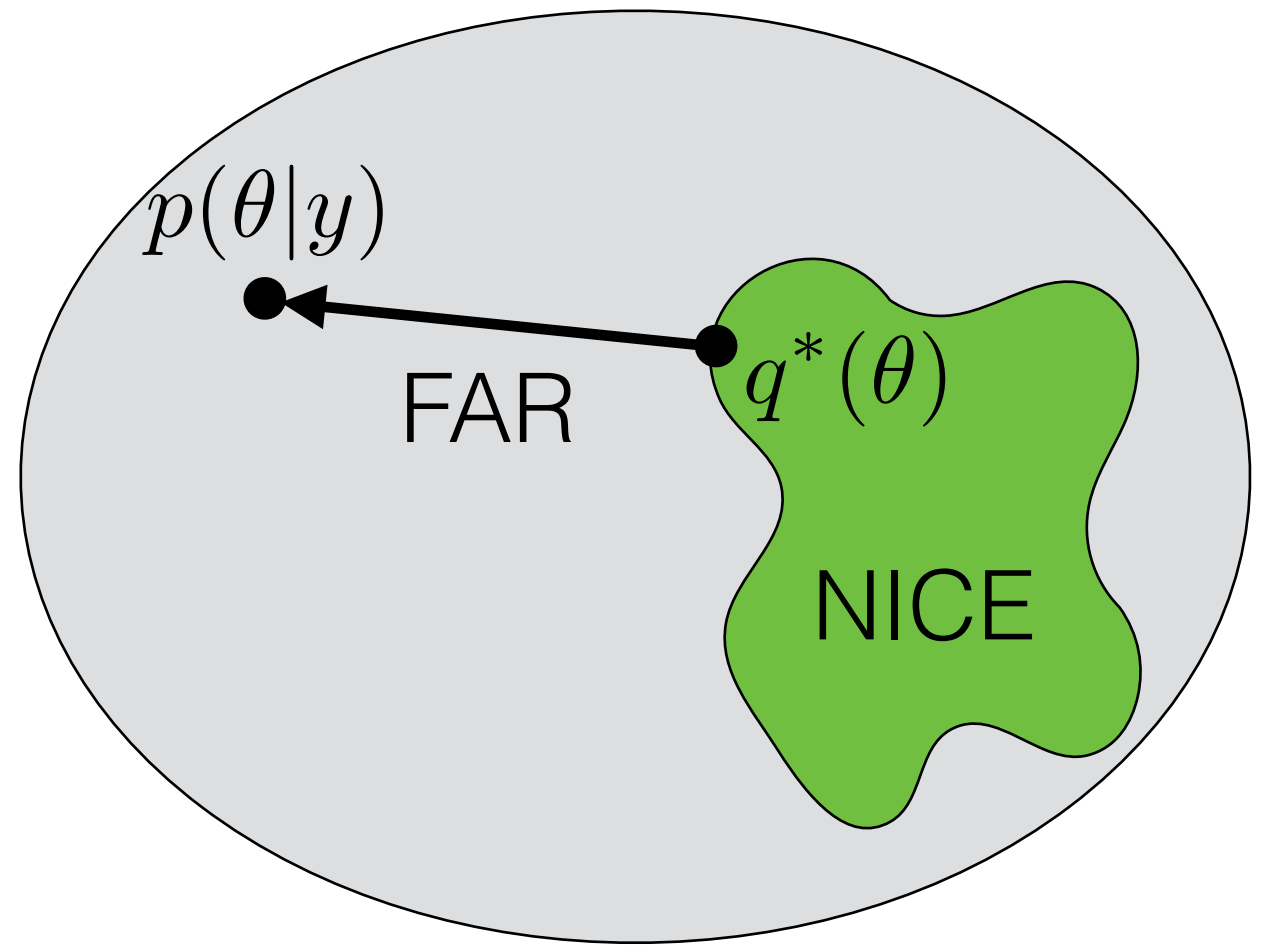
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta$$



Why KL?

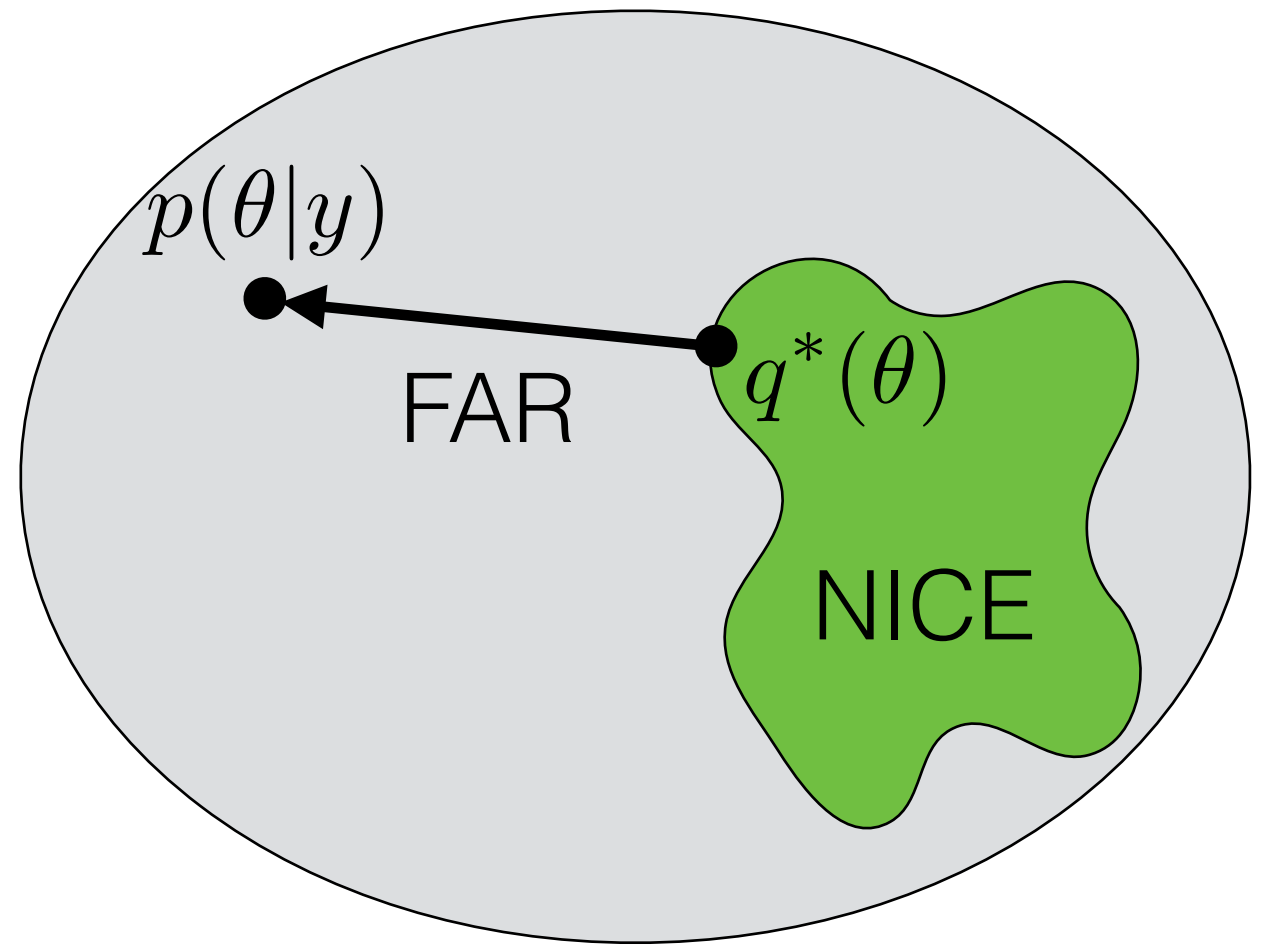
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta$$



Why KL?

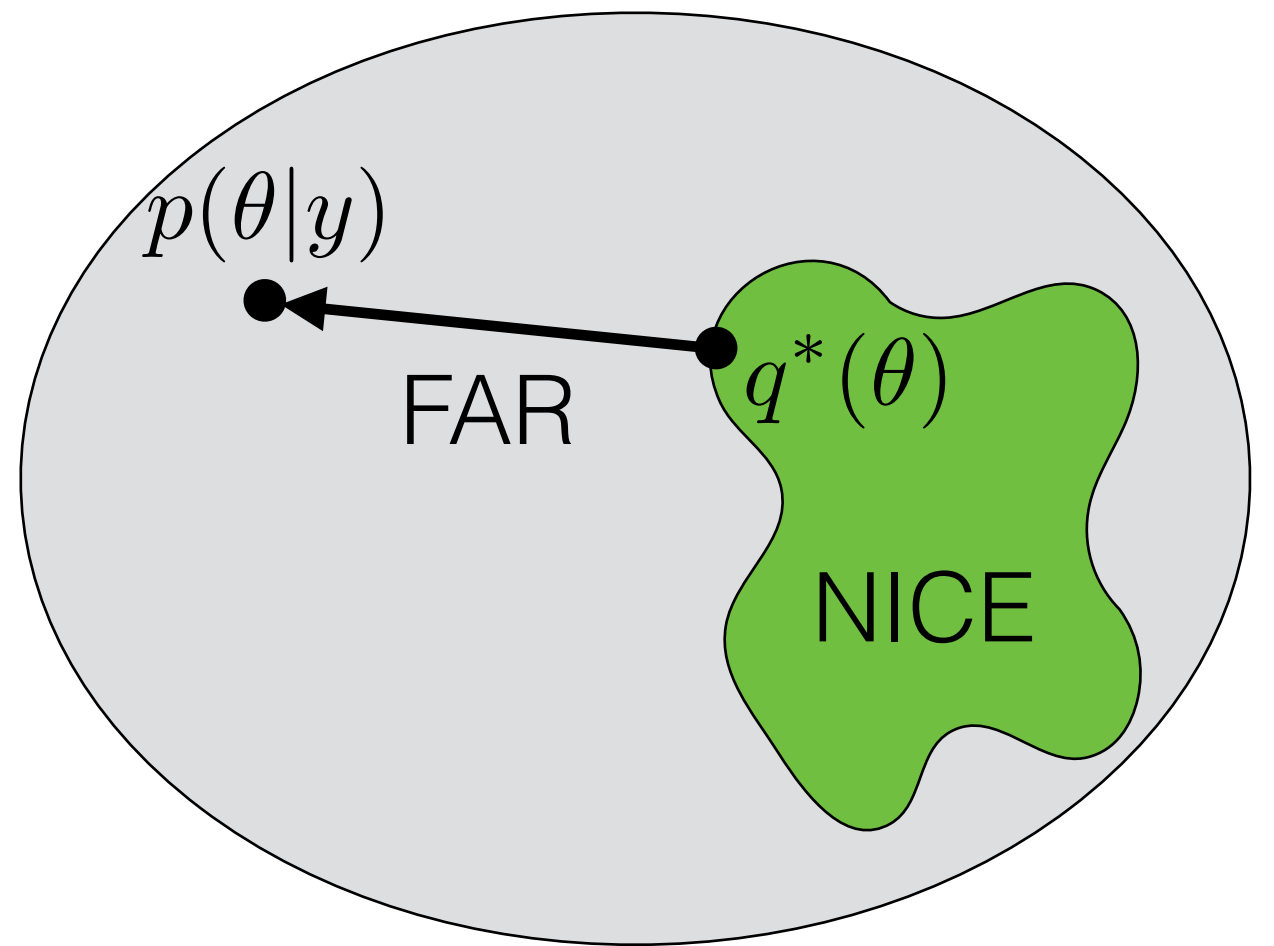
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta$$



Why KL?

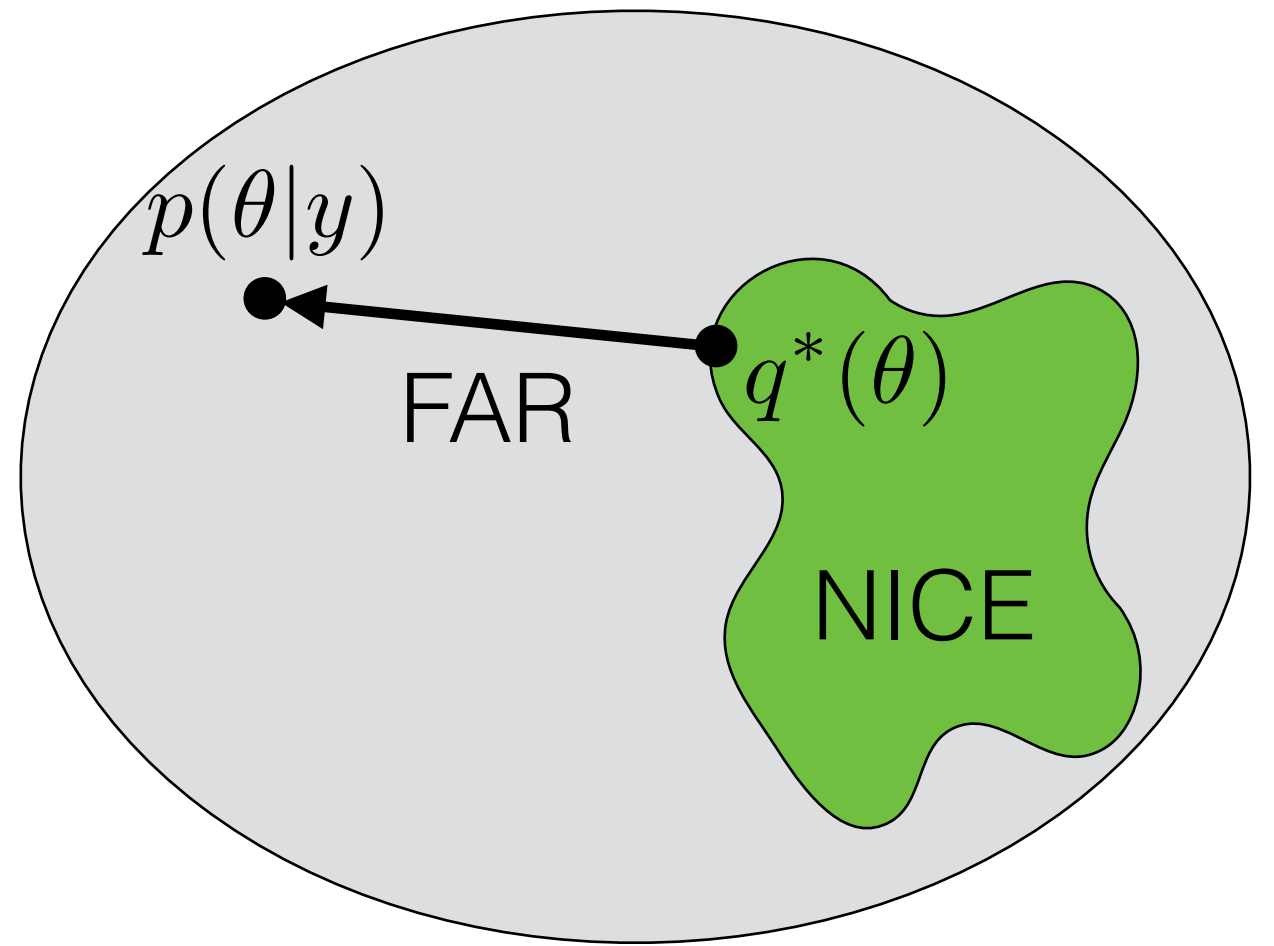
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta$$



Why KL?

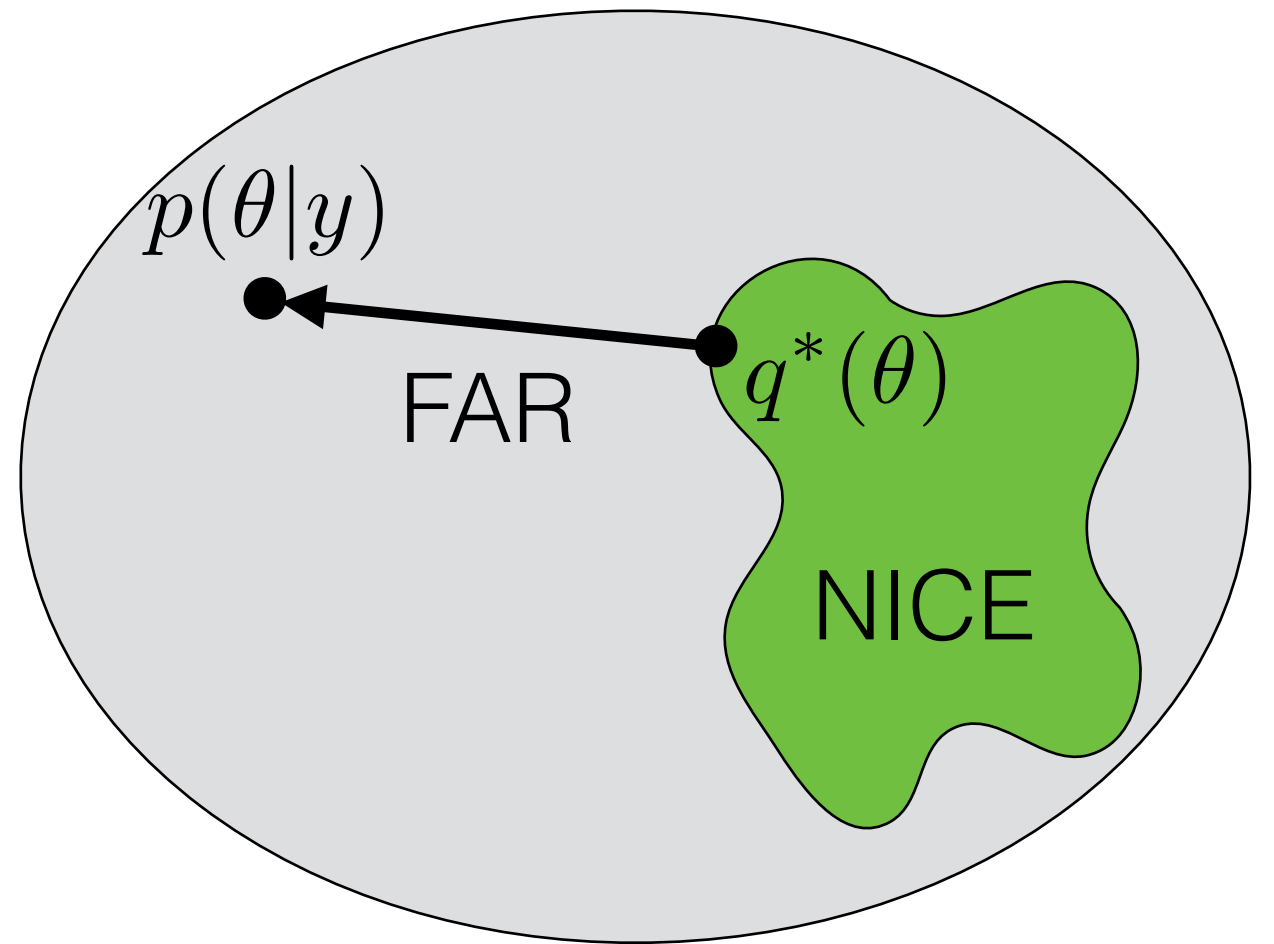
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \int q(\theta) \left[\log p(y) + \log \frac{q(\theta)}{p(\theta, y)} \right] d\theta$$



Why KL?

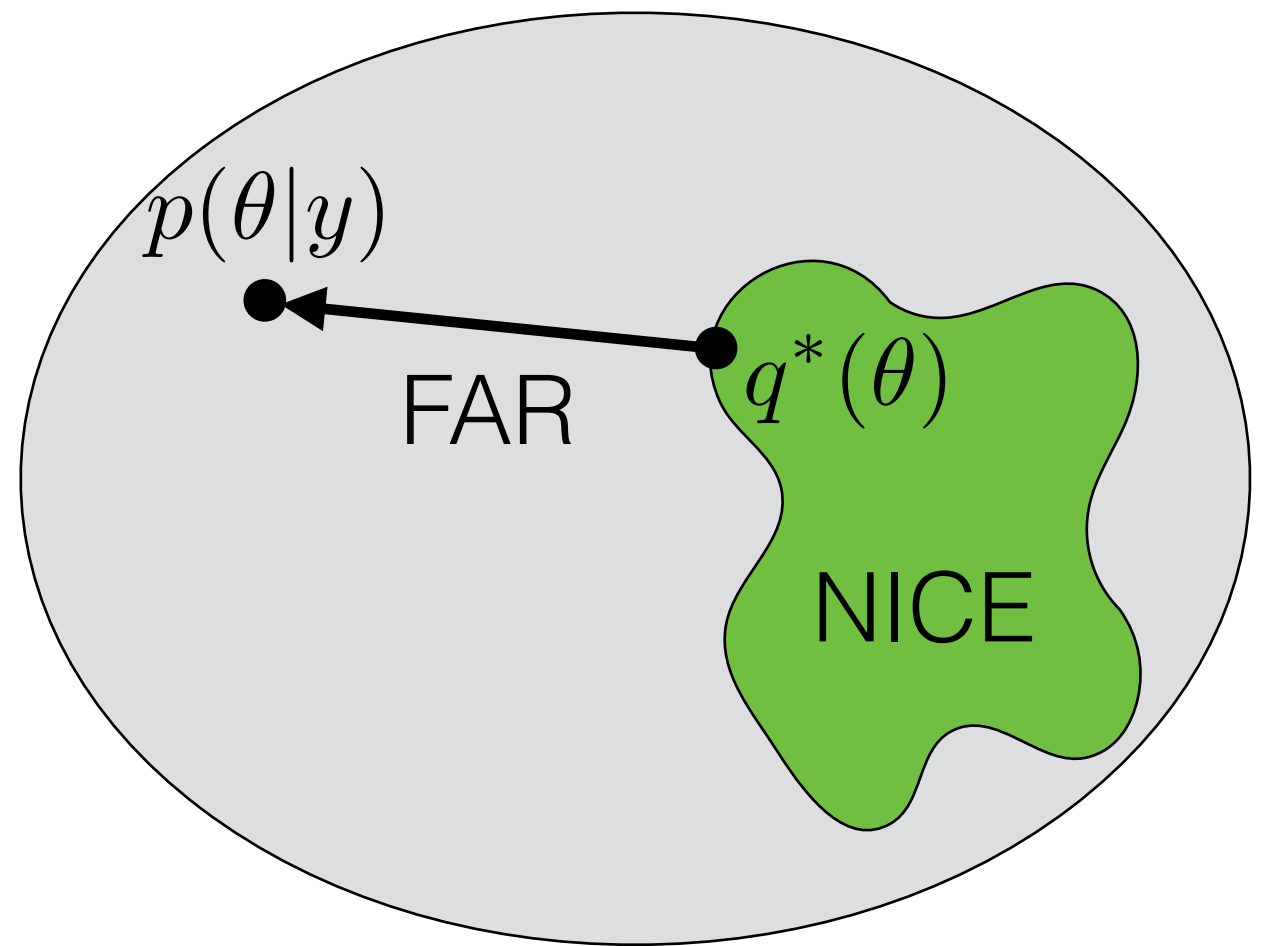
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \int q(\theta) \left[\log p(y) + \log \frac{q(\theta)}{p(\theta, y)} \right] d\theta$$



Why KL?

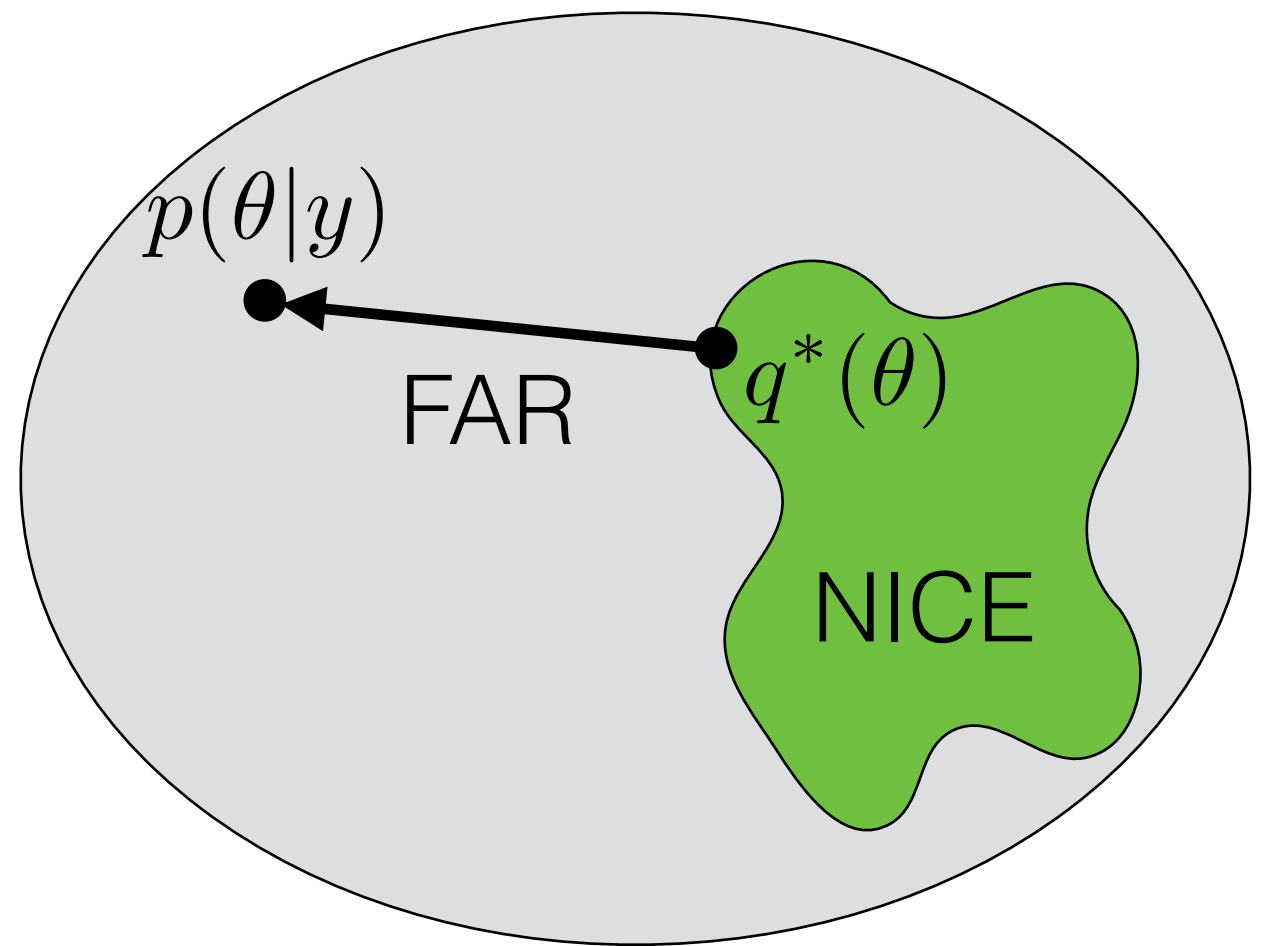
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \int q(\theta) \left[\log p(y) + \log \frac{q(\theta)}{p(\theta, y)} \right] d\theta$$



Why KL?

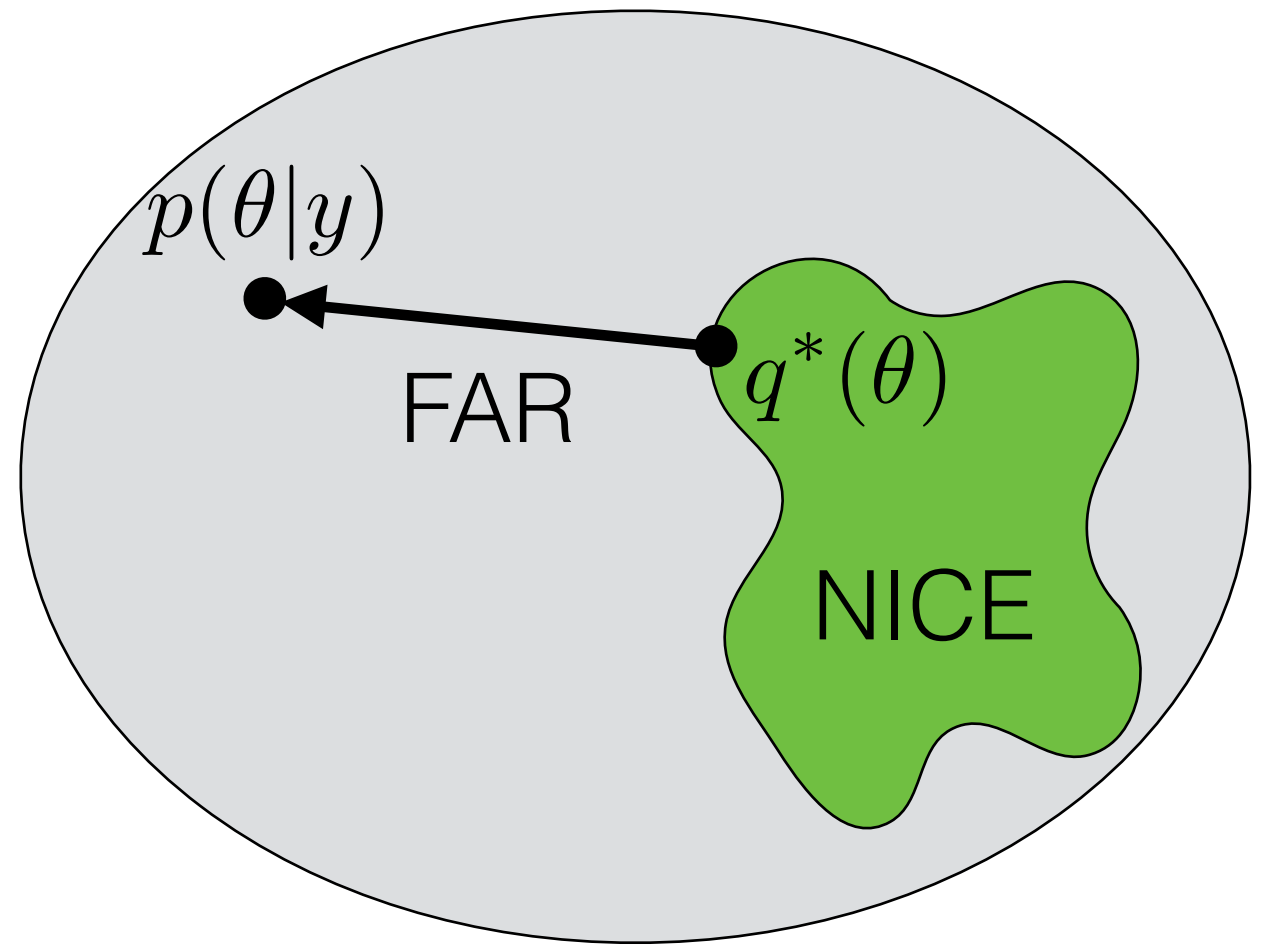
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \int q(\theta) \left[\log p(y) + \log \frac{q(\theta)}{p(\theta, y)} \right] d\theta$$



Why KL?

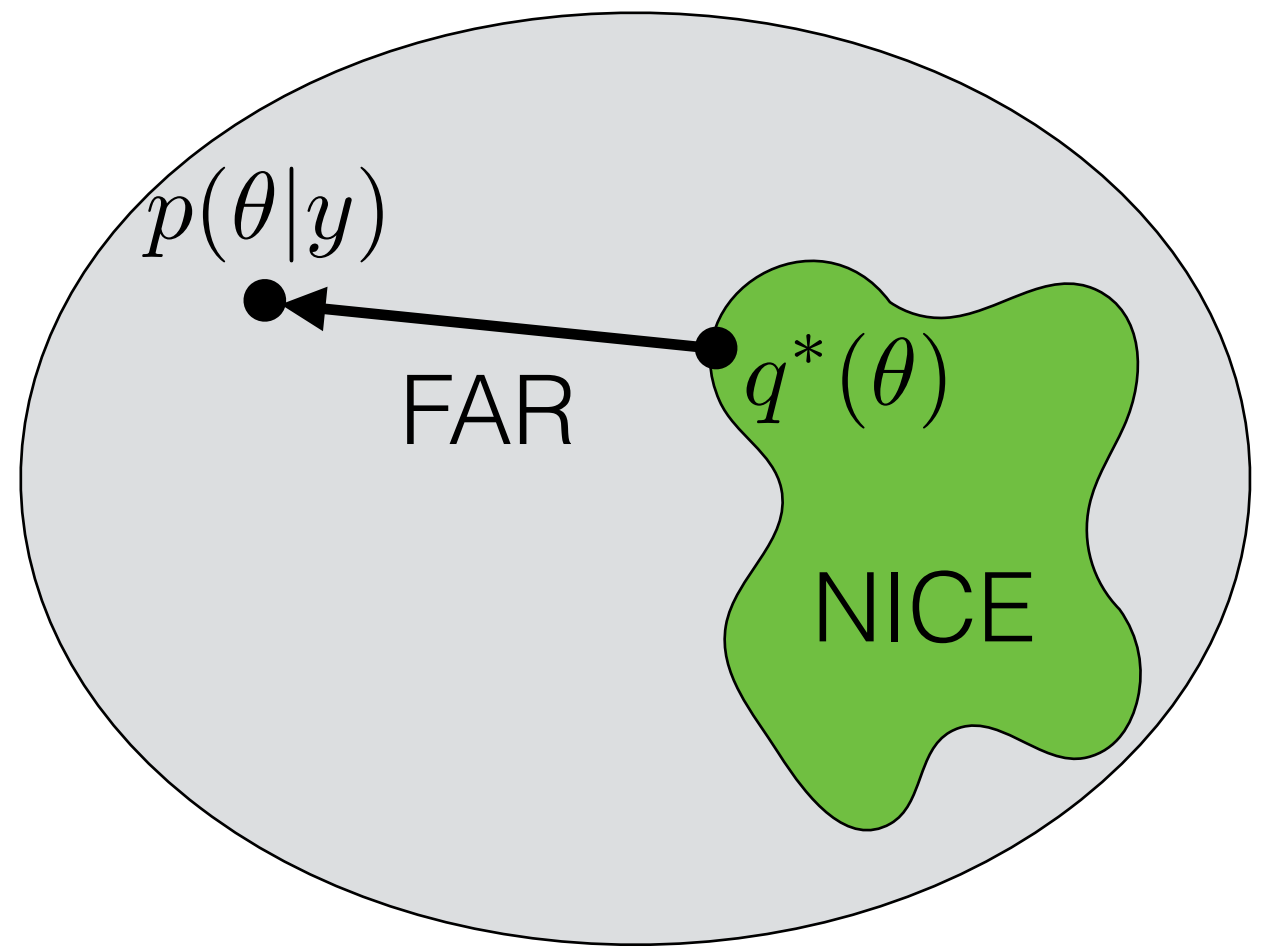
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \int q(\theta) \left[\log p(y) + \log \frac{q(\theta)}{p(\theta, y)} \right] d\theta$$



Why KL?

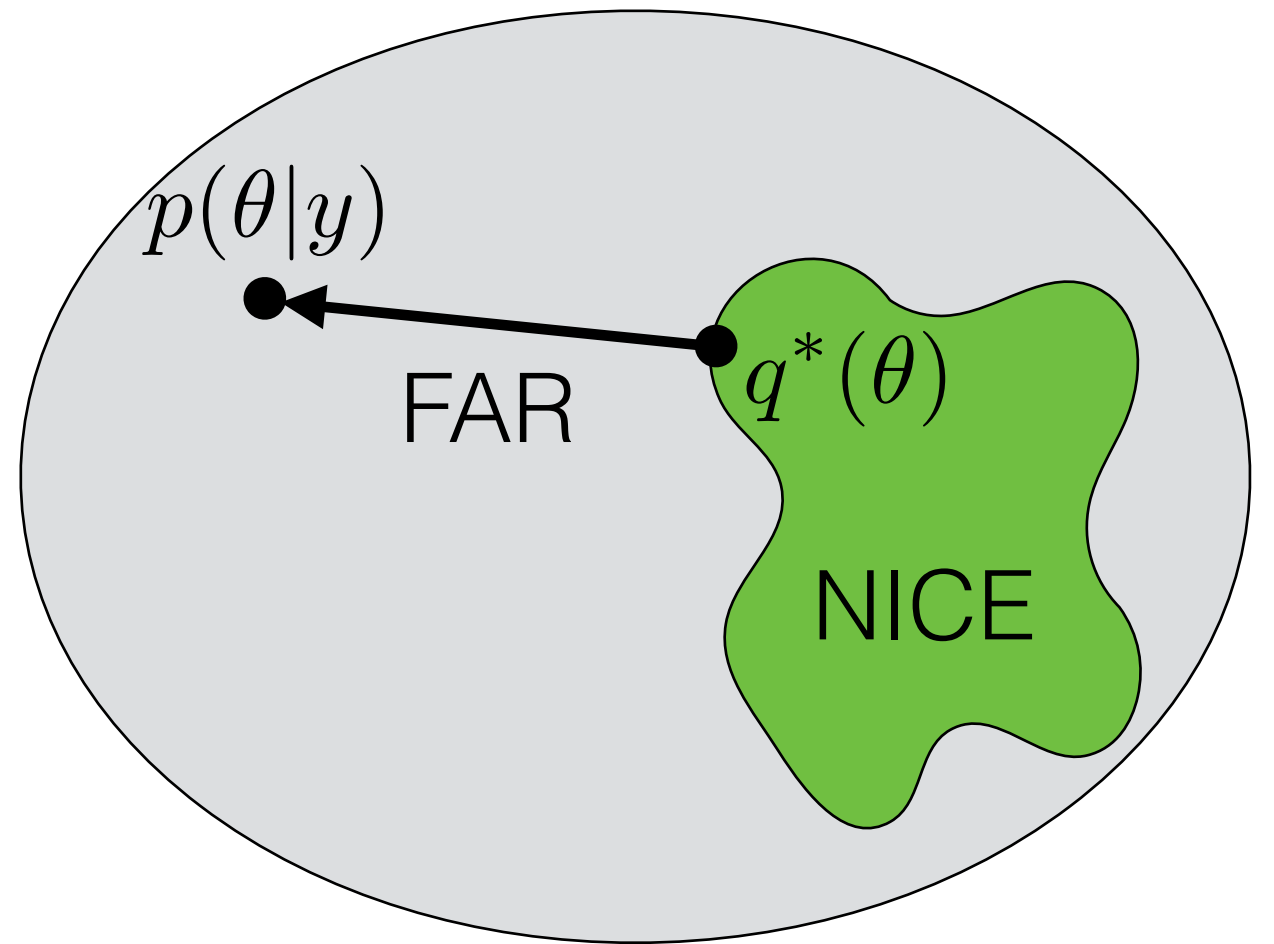
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \int q(\theta) \left[\log p(y) + \log \frac{q(\theta)}{p(\theta, y)} \right] d\theta$$



Why KL?

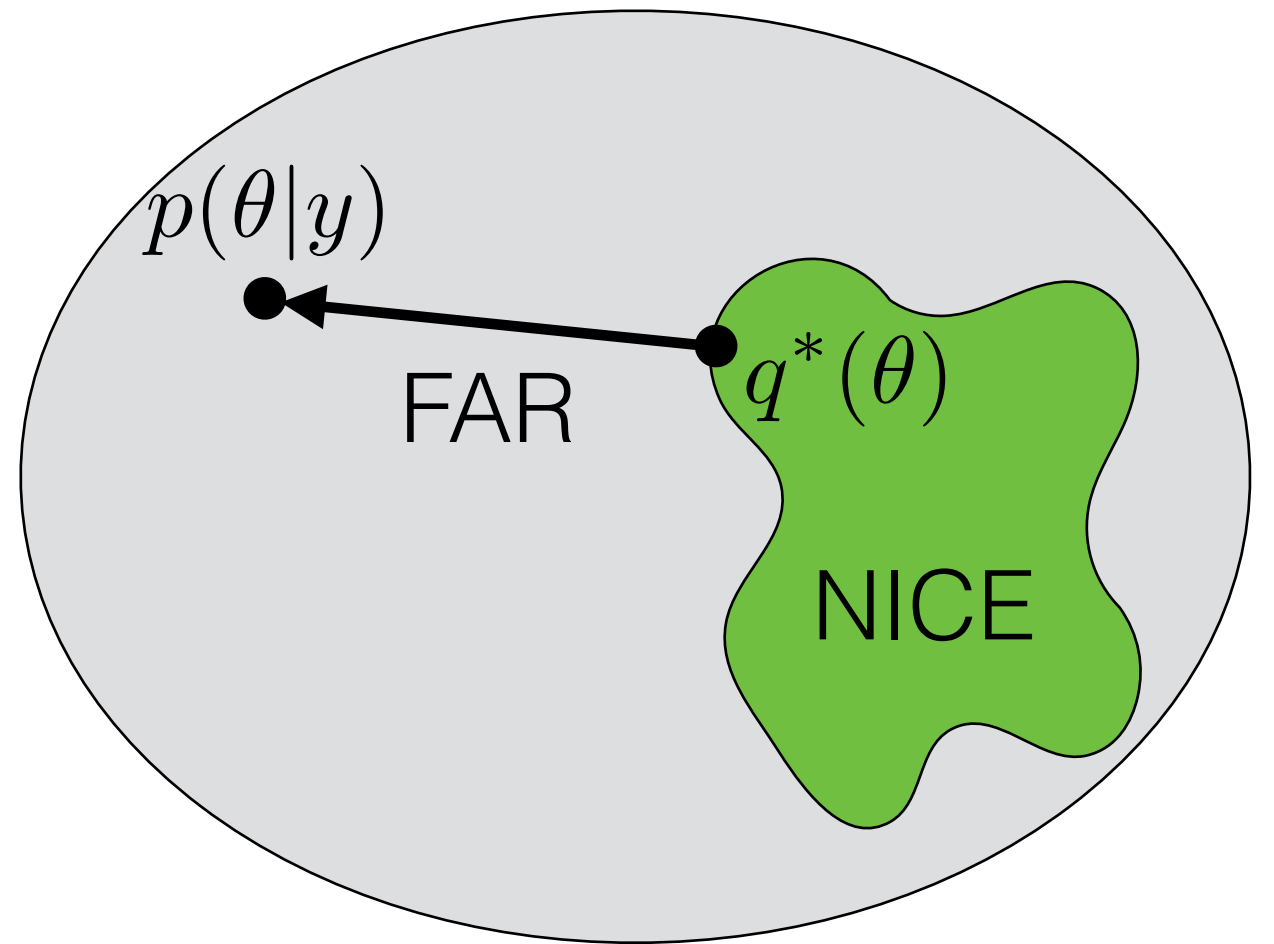
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL} (q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL} (q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \int \boxed{q(\theta) \left[\log p(y) + \log \frac{q(\theta)}{p(\theta, y)} \right]} d\theta$$



Why KL?

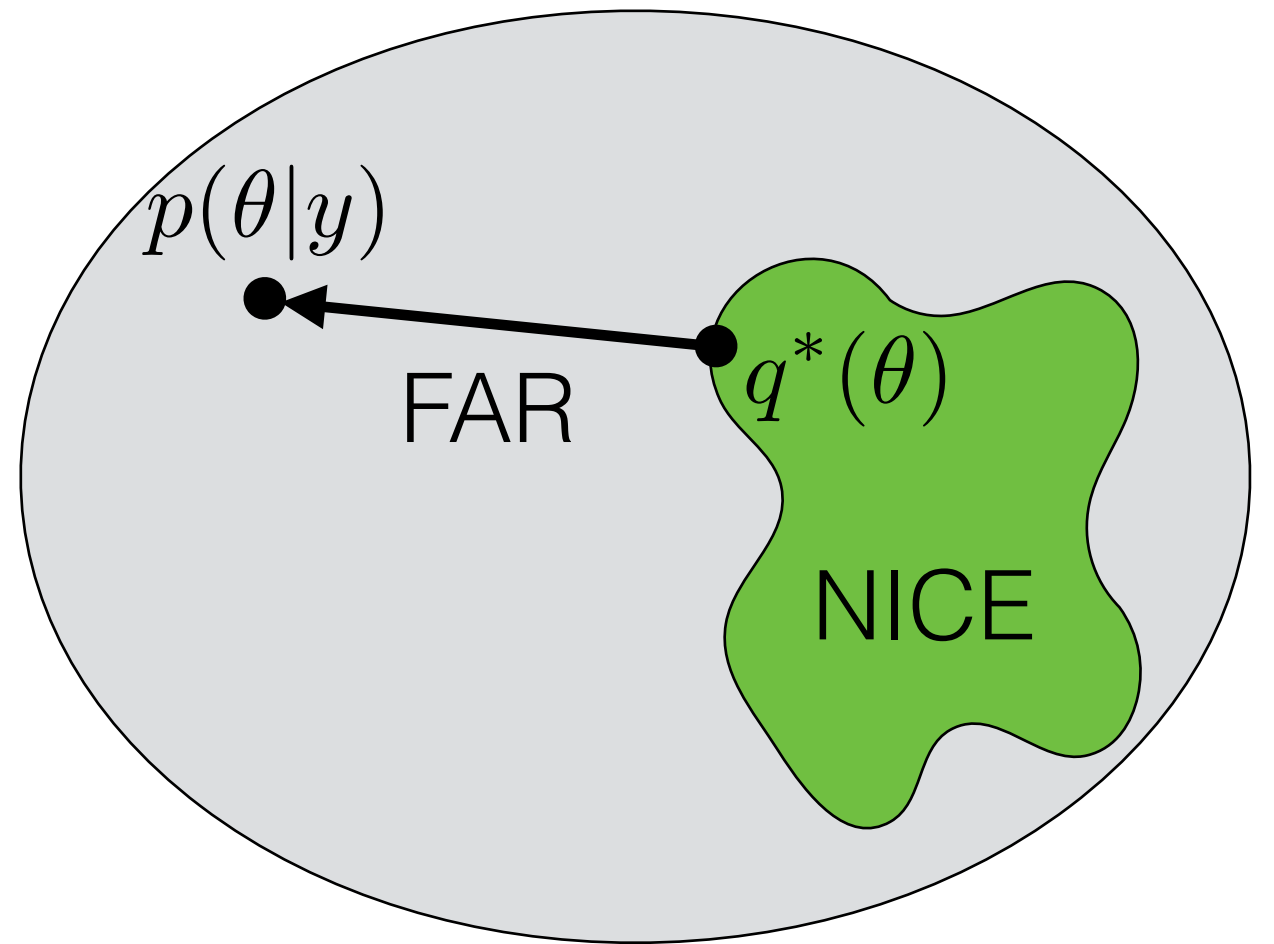
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) + \int q(\theta) \log \frac{q(\theta)}{p(\theta, y)} d\theta$$



Why KL?

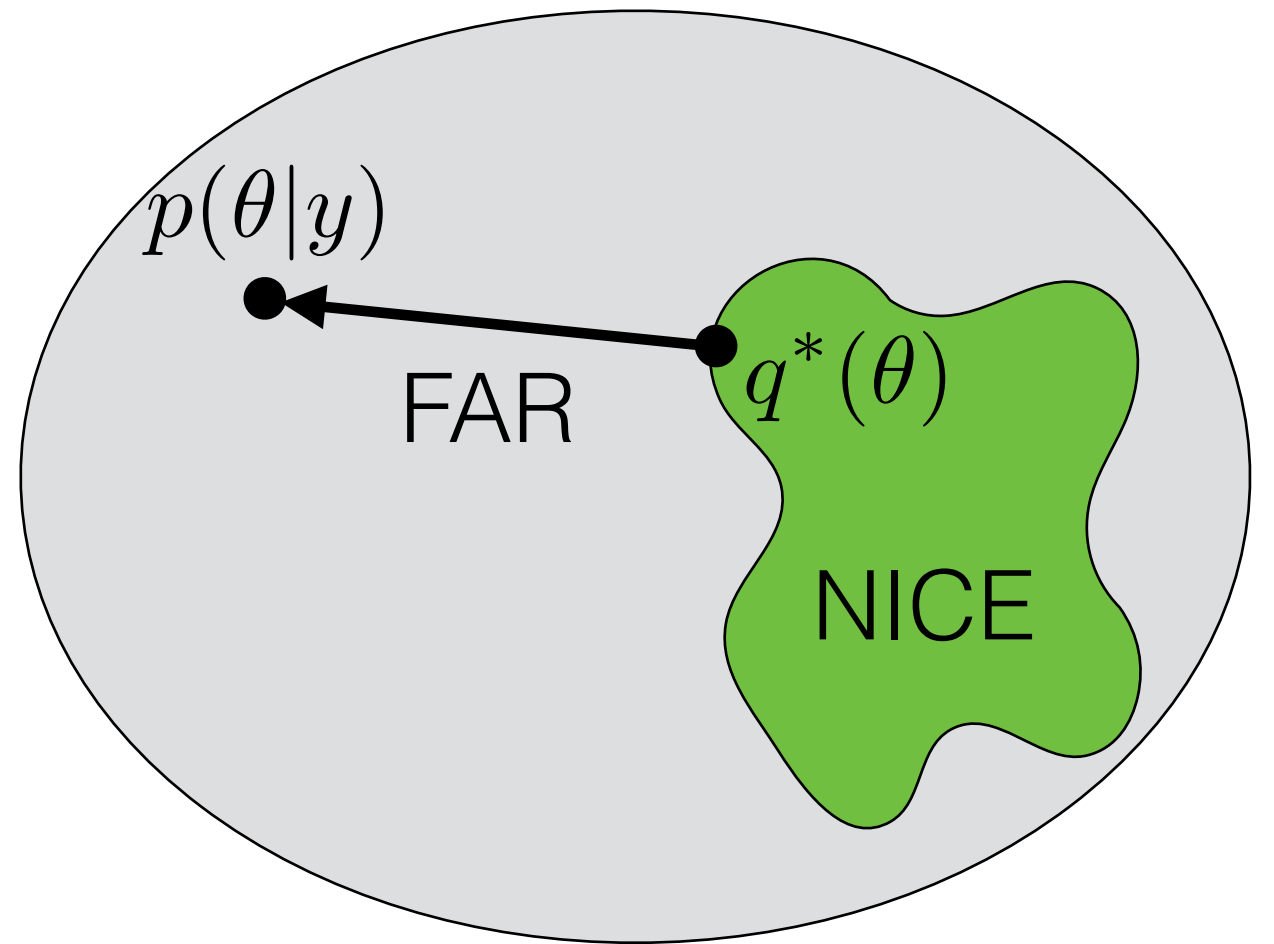
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) \boxed{+} \int q(\theta) \log \boxed{\frac{q(\theta)}{p(\theta, y)}} d\theta$$



Why KL?

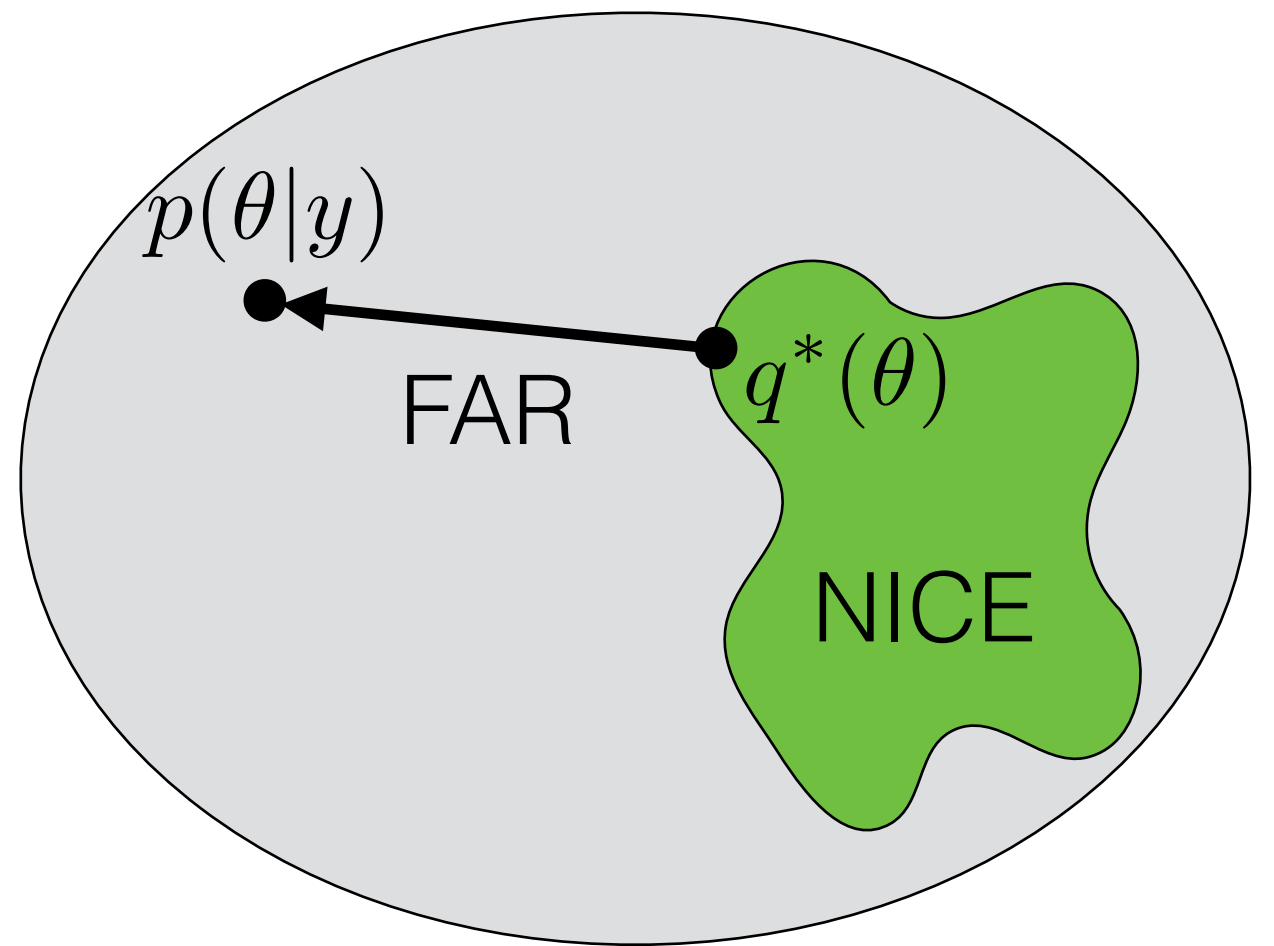
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) \text{ — } \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$



Why KL?

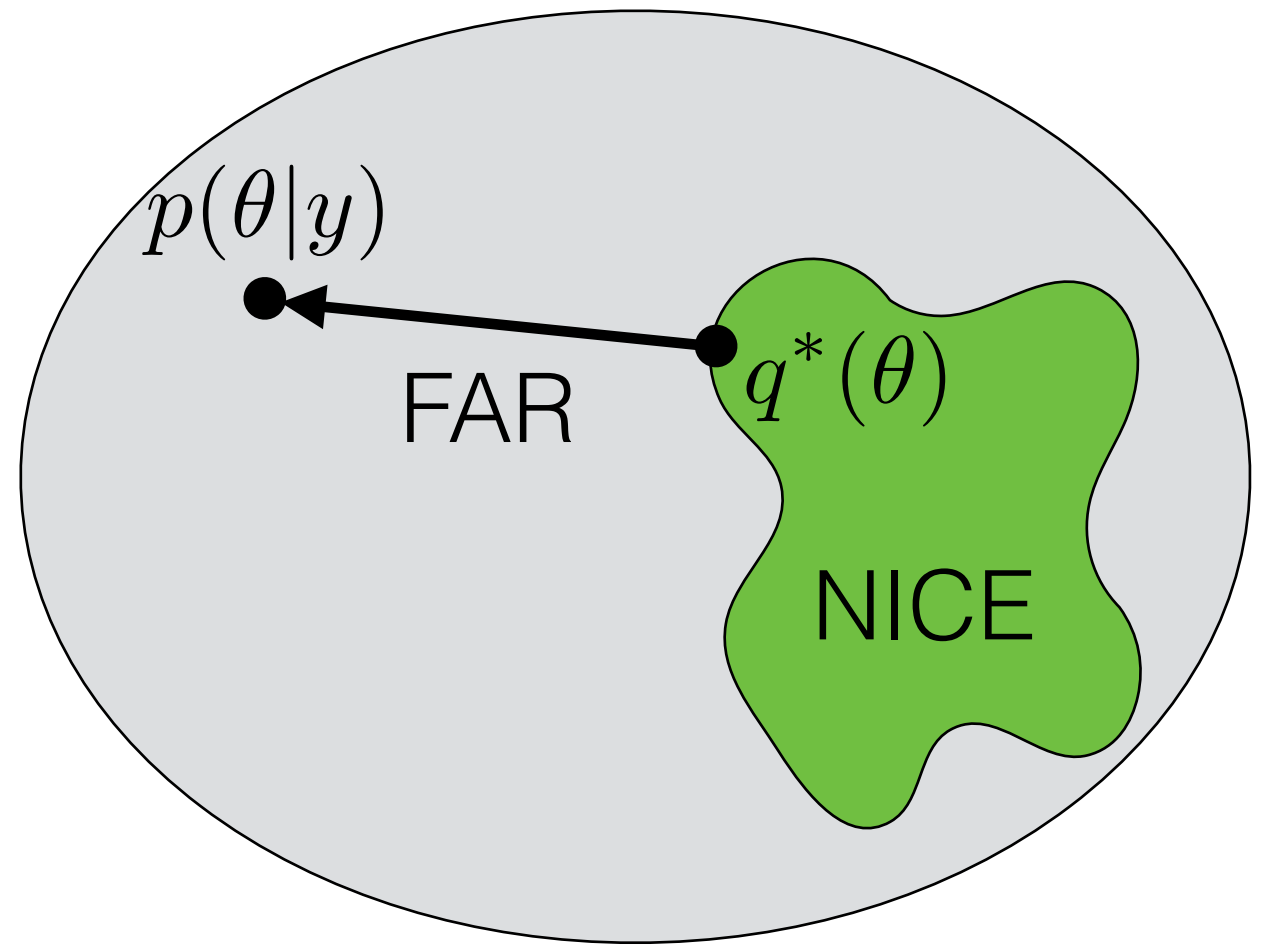
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$



Why KL?

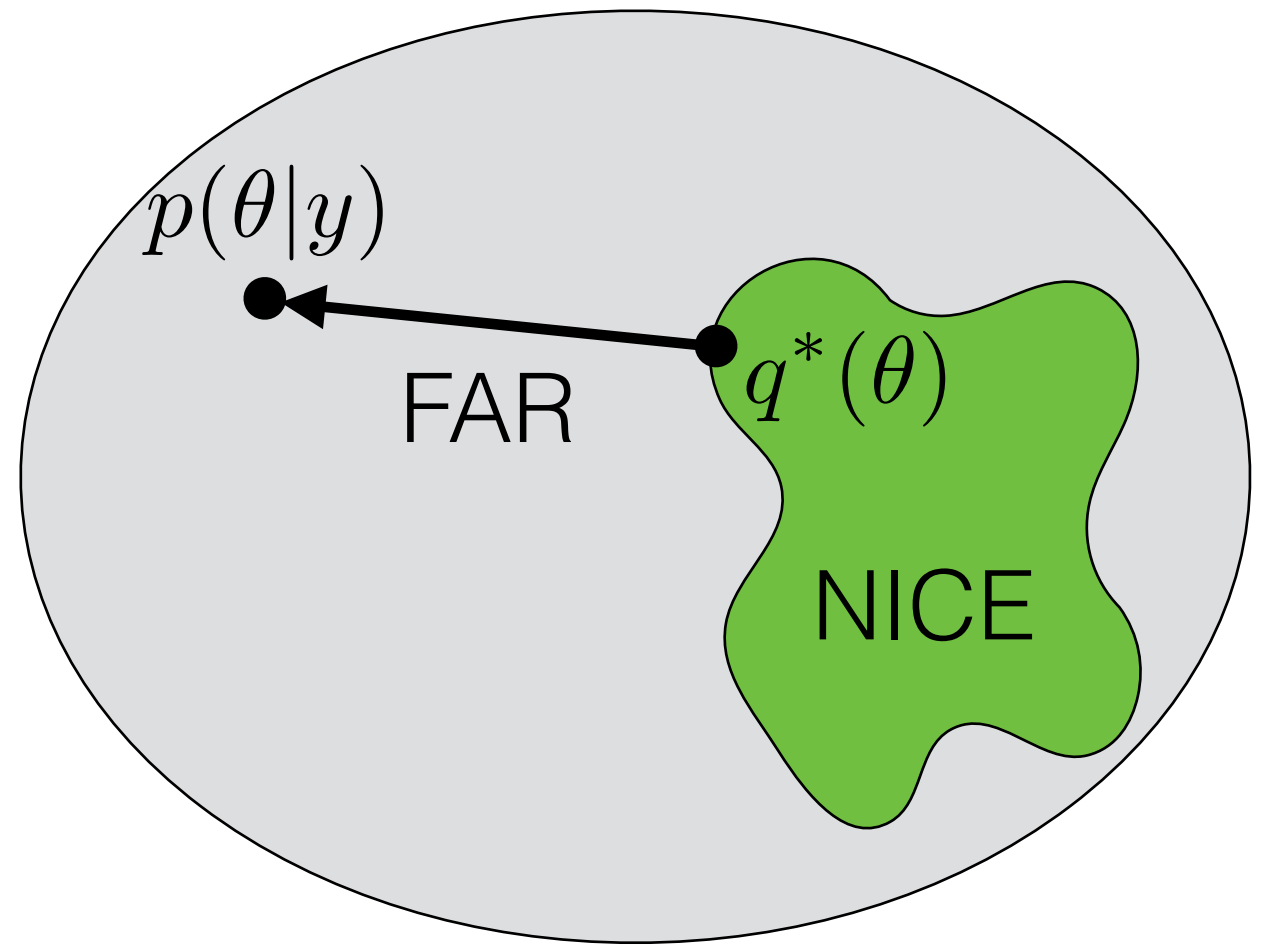
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$



Why KL?

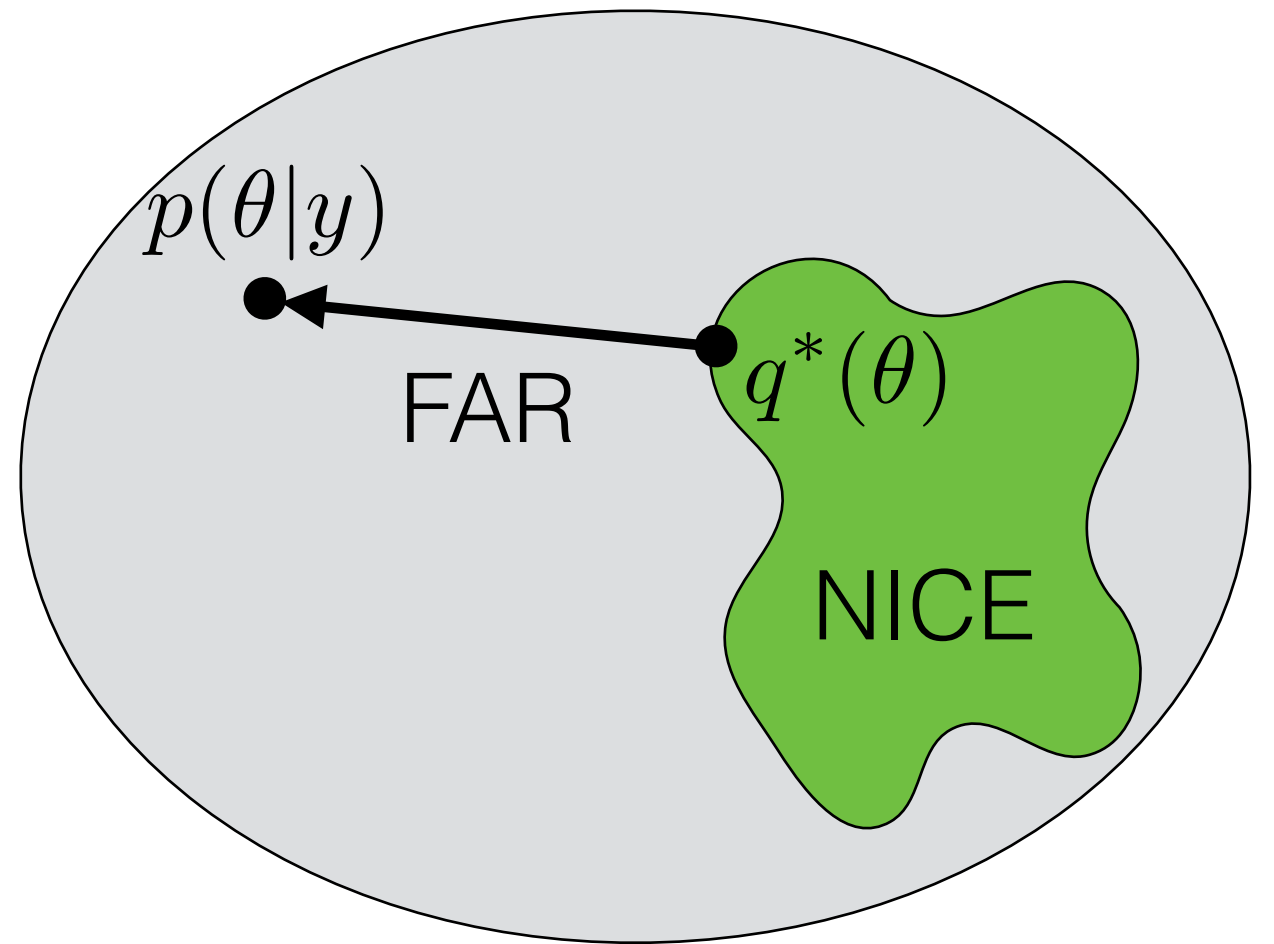
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$



Why KL?

- Variational Bayes

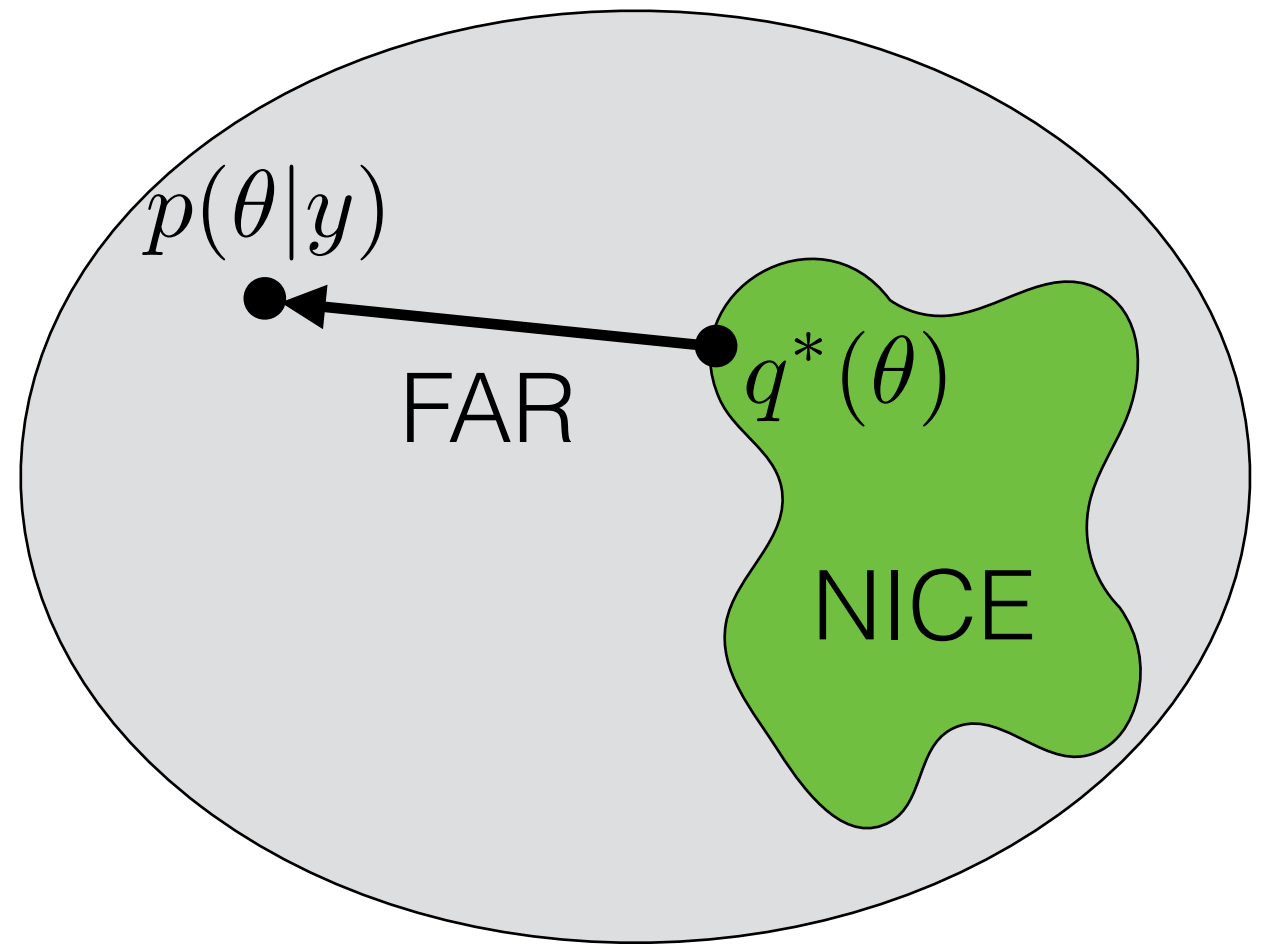
$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$

“Evidence lower bound” (ELBO)



Why KL?

- Variational Bayes

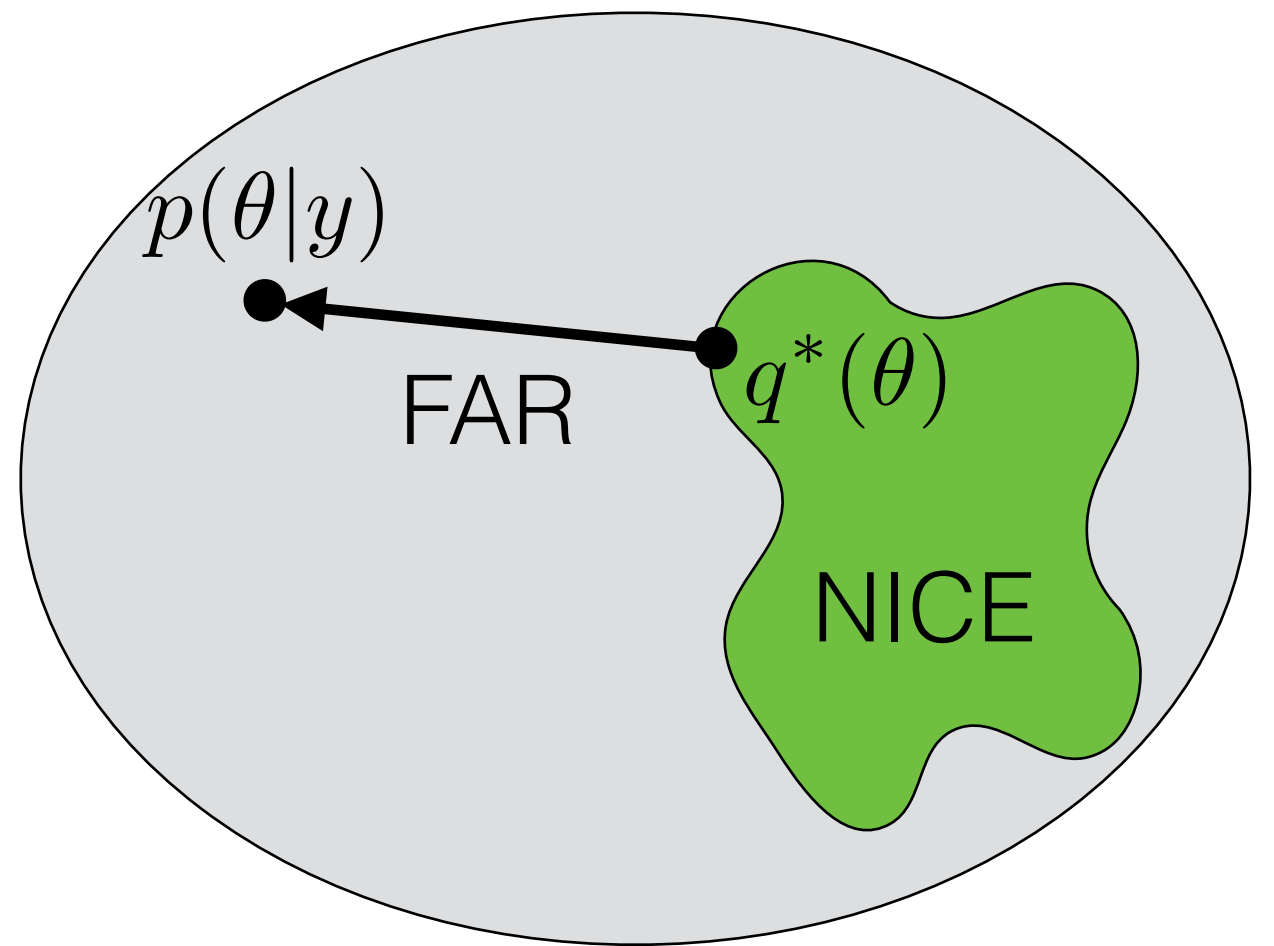
$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$

“Evidence lower bound” (ELBO)



Why KL?

- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

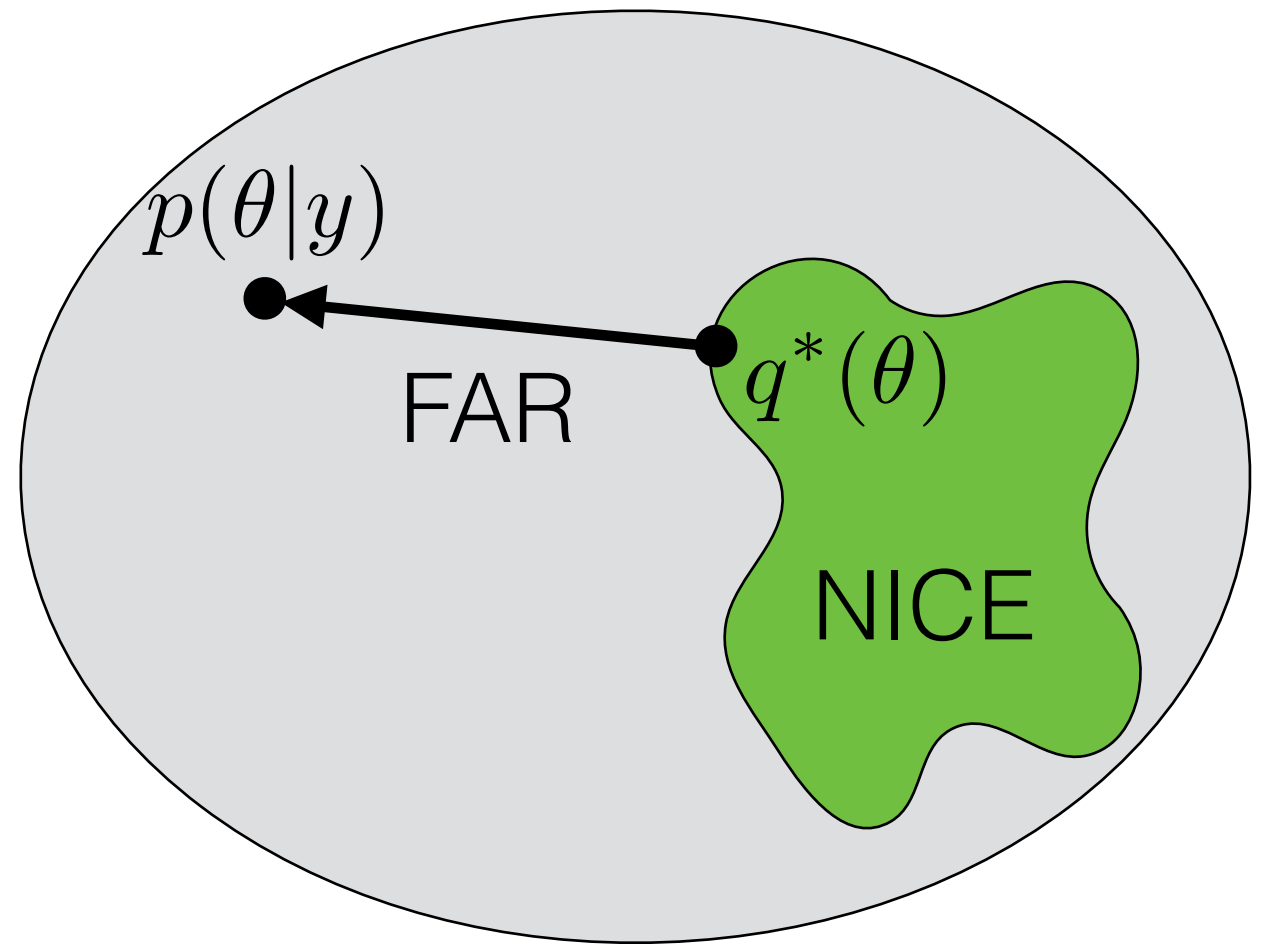
$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$

- Exercise: Show $\operatorname{KL} \geq 0$ [Bishop 2006, Sec 1.6.1]

“Evidence lower bound” (ELBO)



Why KL?

- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

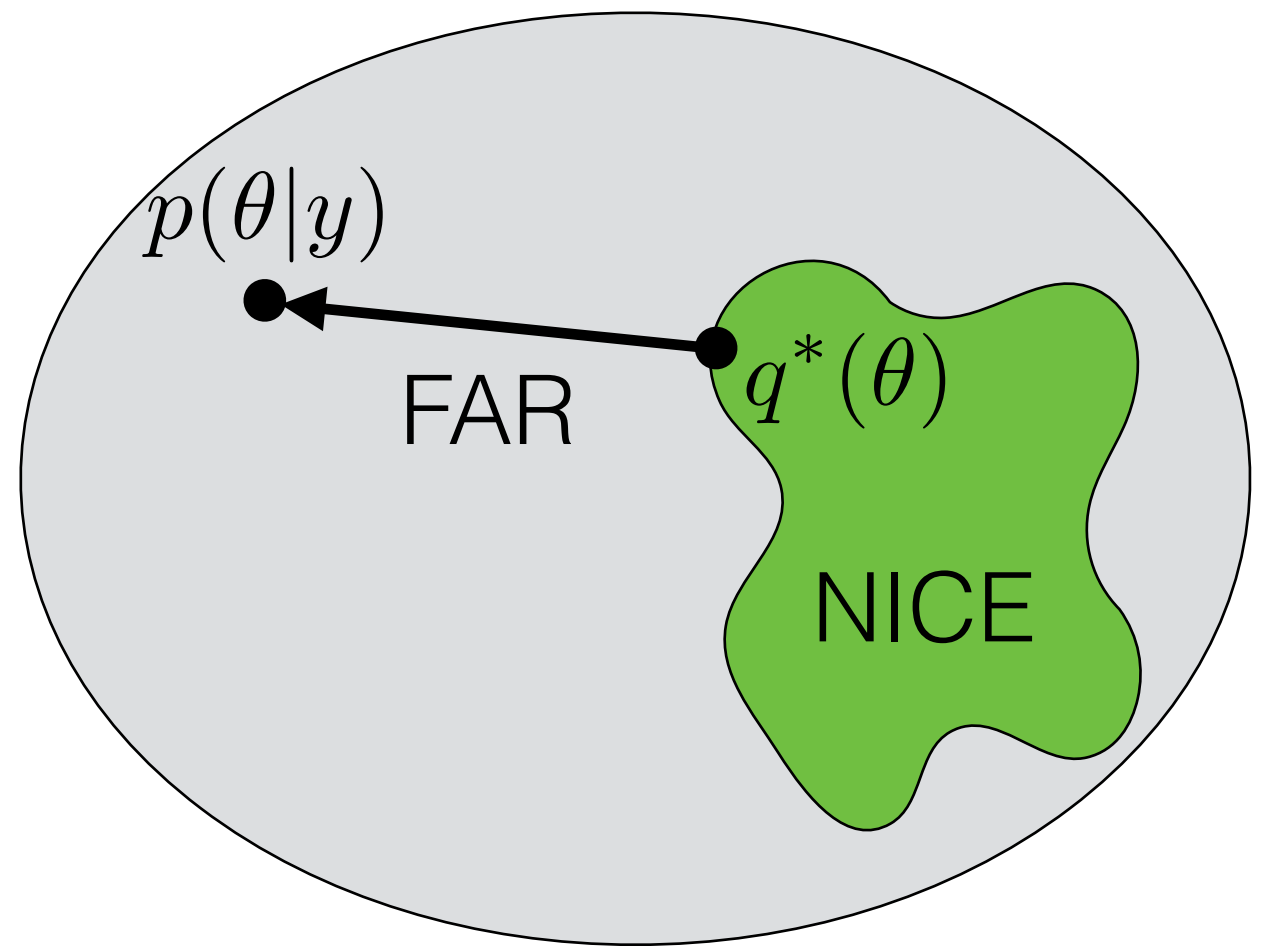
$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$

- Exercise: Show $\operatorname{KL} \geq 0$ [Bishop 2006, Sec 1.6.1]

- $\operatorname{KL} \geq 0 \Rightarrow \log p(y) \geq \operatorname{ELBO}$



“Evidence lower bound” (ELBO)

Why KL?

- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

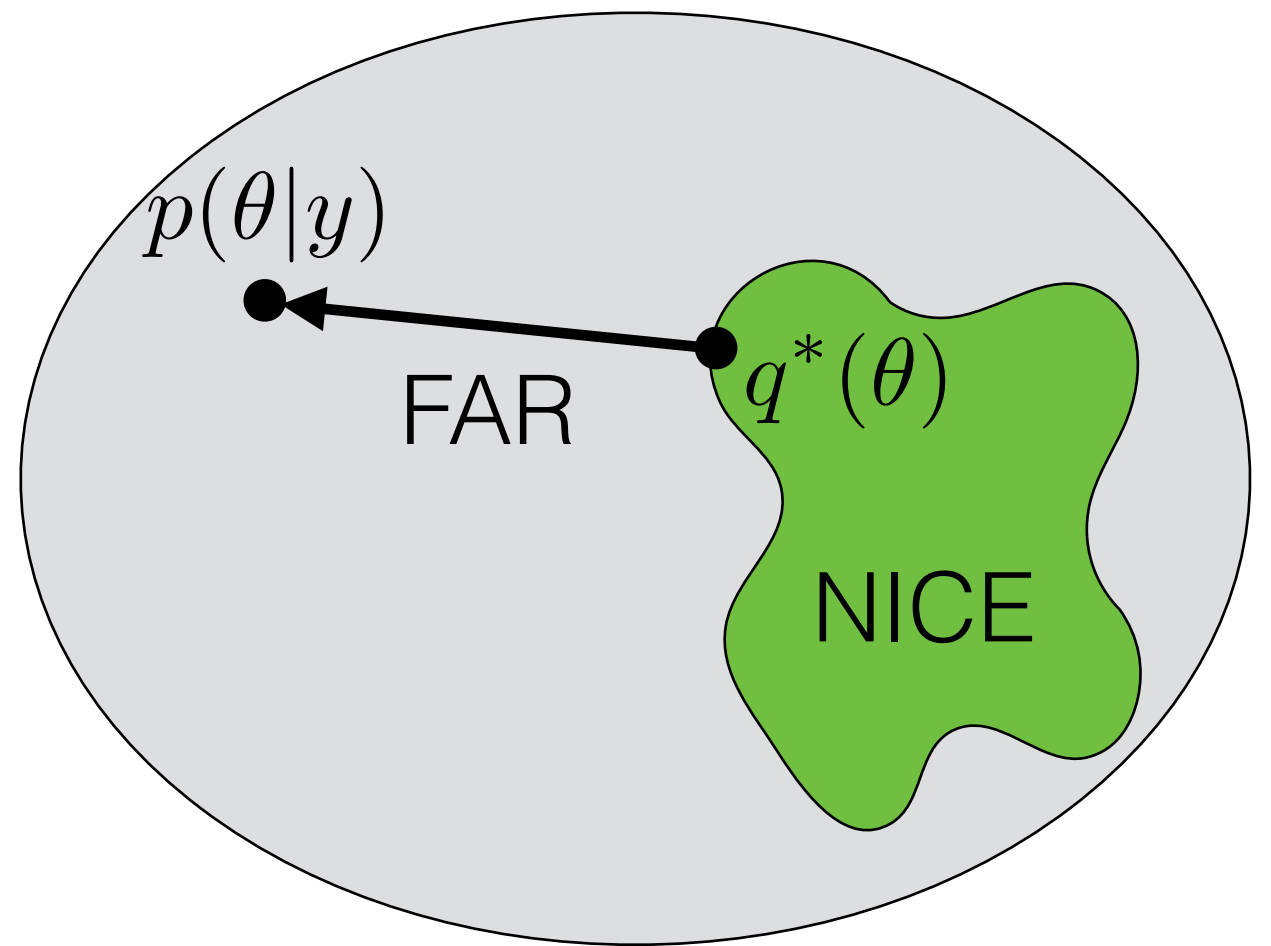
$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$

- Exercise: Show $\operatorname{KL} \geq 0$ [Bishop 2006, Sec 1.6.1]

- $\operatorname{KL} \geq 0 \Rightarrow \log p(y) \geq \operatorname{ELBO}$

- $q^* = \operatorname{argmax}_{q \in Q} \operatorname{ELBO}(q)$



“Evidence lower bound” (ELBO)

Why KL?

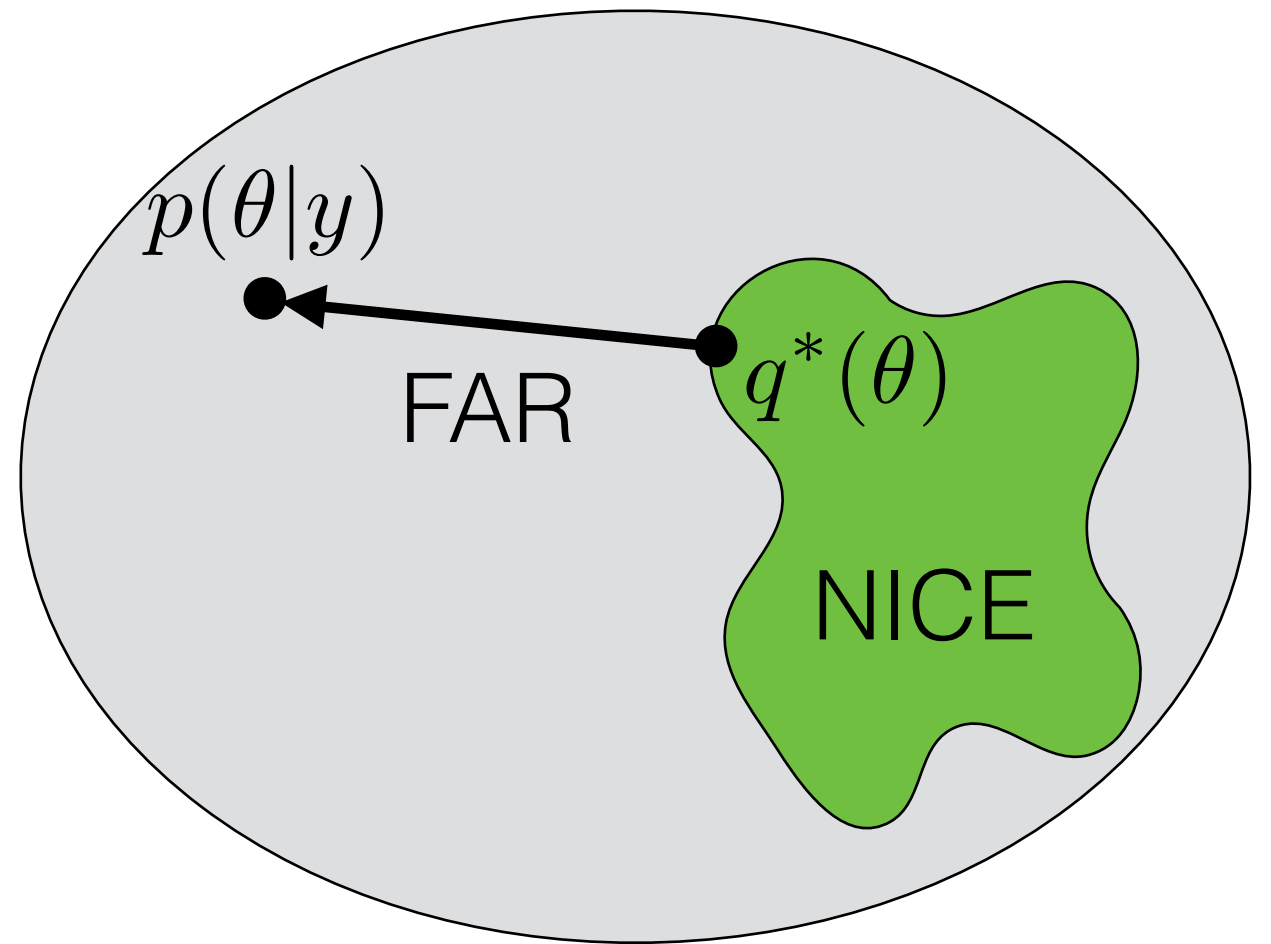
- Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$\operatorname{KL}(q(\cdot) || p(\cdot|y))$$

$$:= \int q(\theta) \log \frac{q(\theta)}{p(\theta|y)} d\theta$$

$$= \int q(\theta) \log \frac{q(\theta)p(y)}{p(\theta, y)} d\theta = \log p(y) - \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$



- Exercise: Show $\operatorname{KL} \geq 0$ [Bishop 2006, Sec 1.6.1]

- $\operatorname{KL} \geq 0 \Rightarrow \log p(y) \geq \operatorname{ELBO}$

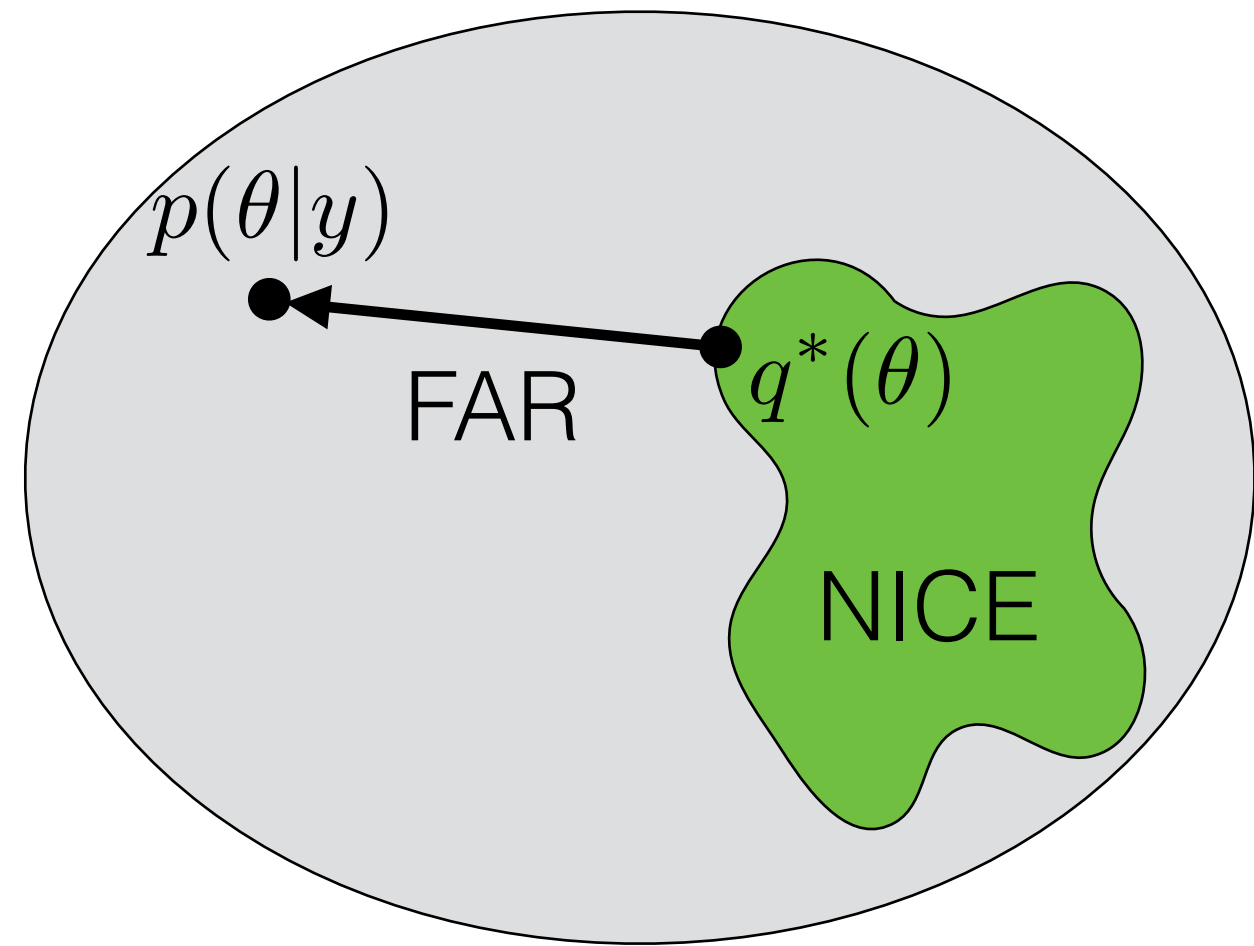
- $q^* = \operatorname{argmax}_{q \in Q} \operatorname{ELBO}(q)$

- Why KL (in this direction)?

“Evidence lower bound” (ELBO)

Variational Bayes

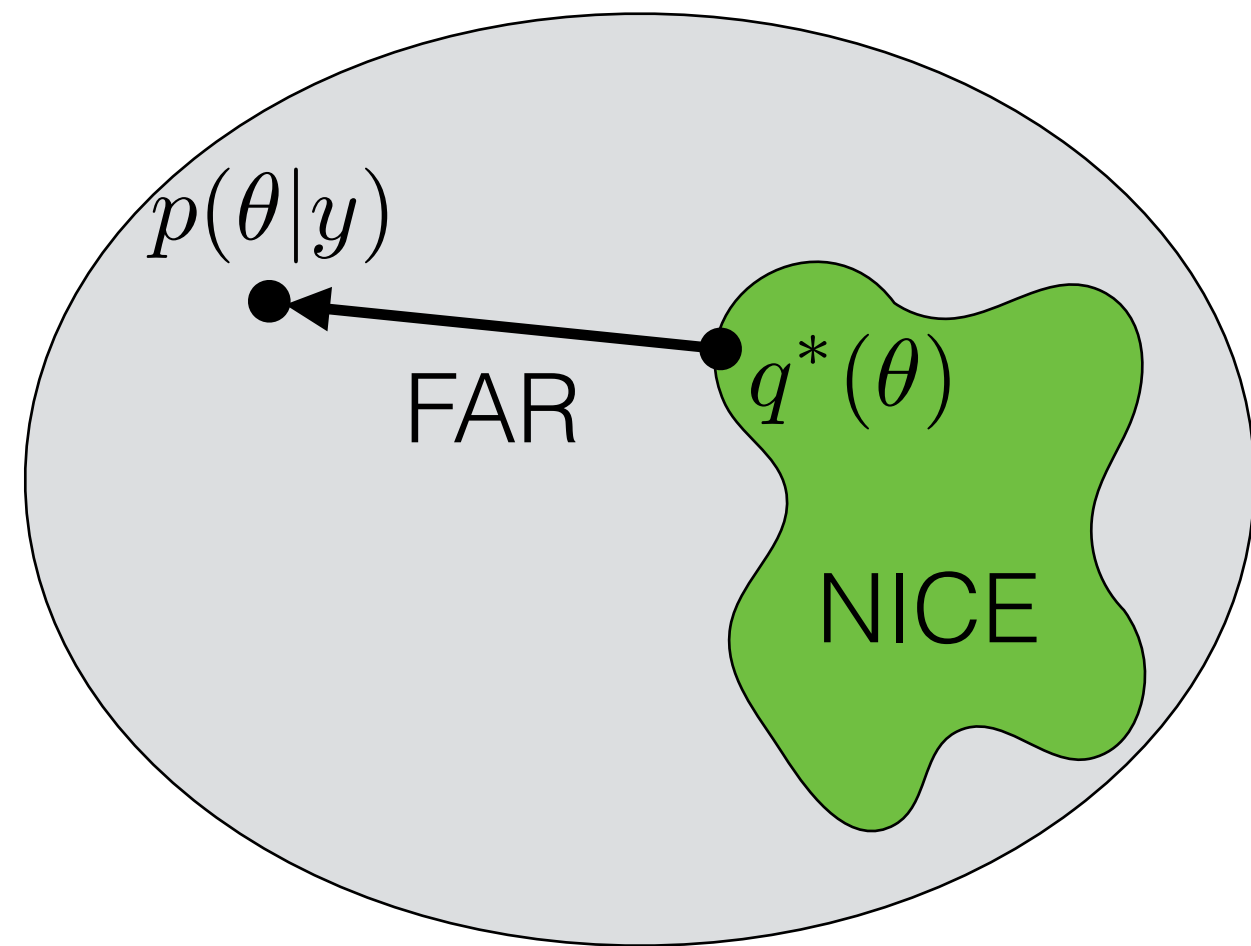
$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL} (q(\cdot) || p(\cdot|y))$$



Variational Bayes

$$q^* = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

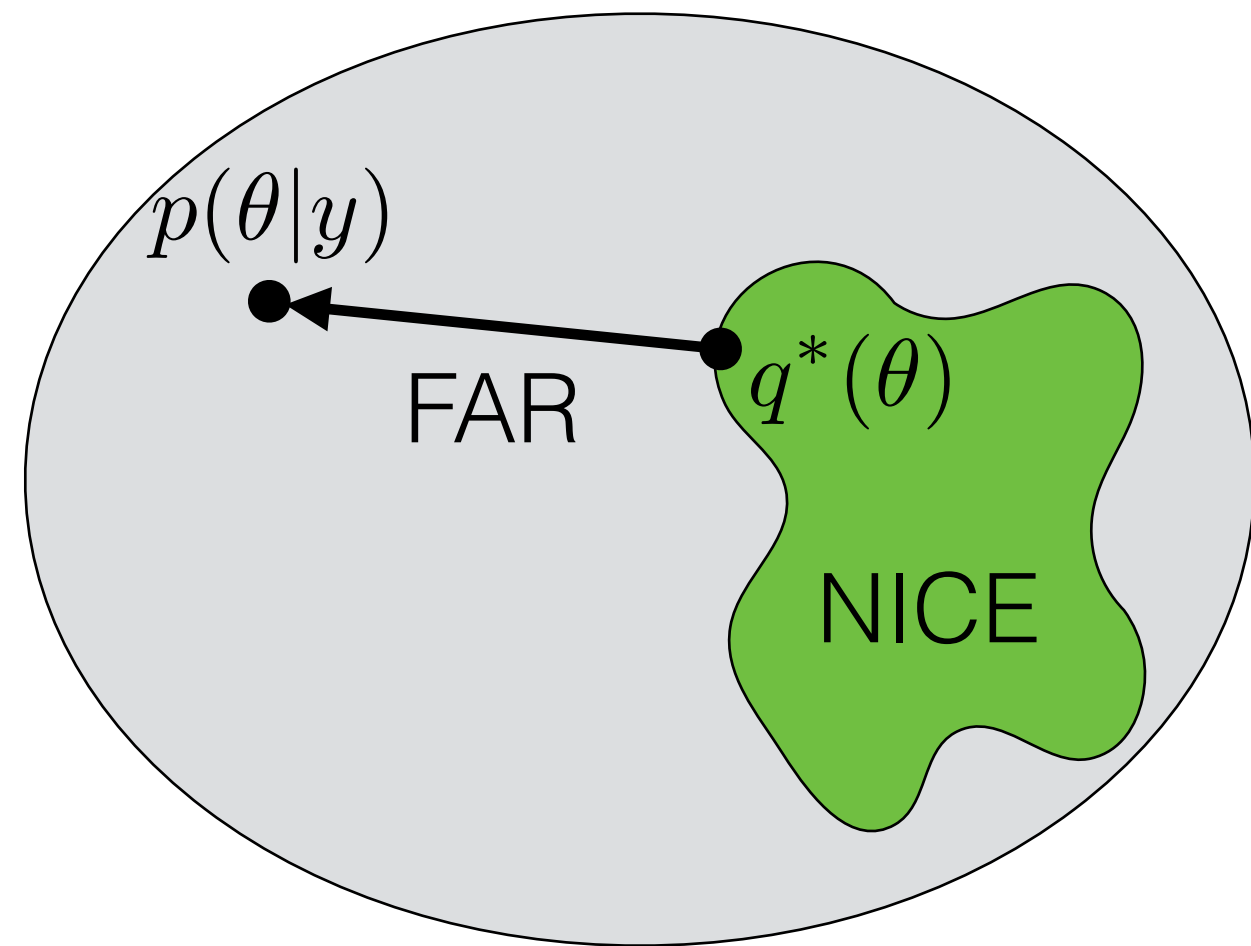
Choose “NICE” distributions



Variational Bayes

$$q^* = \operatorname{argmin}_{q \in \mathcal{Q}} \text{KL}(q(\cdot) || p(\cdot|y))$$

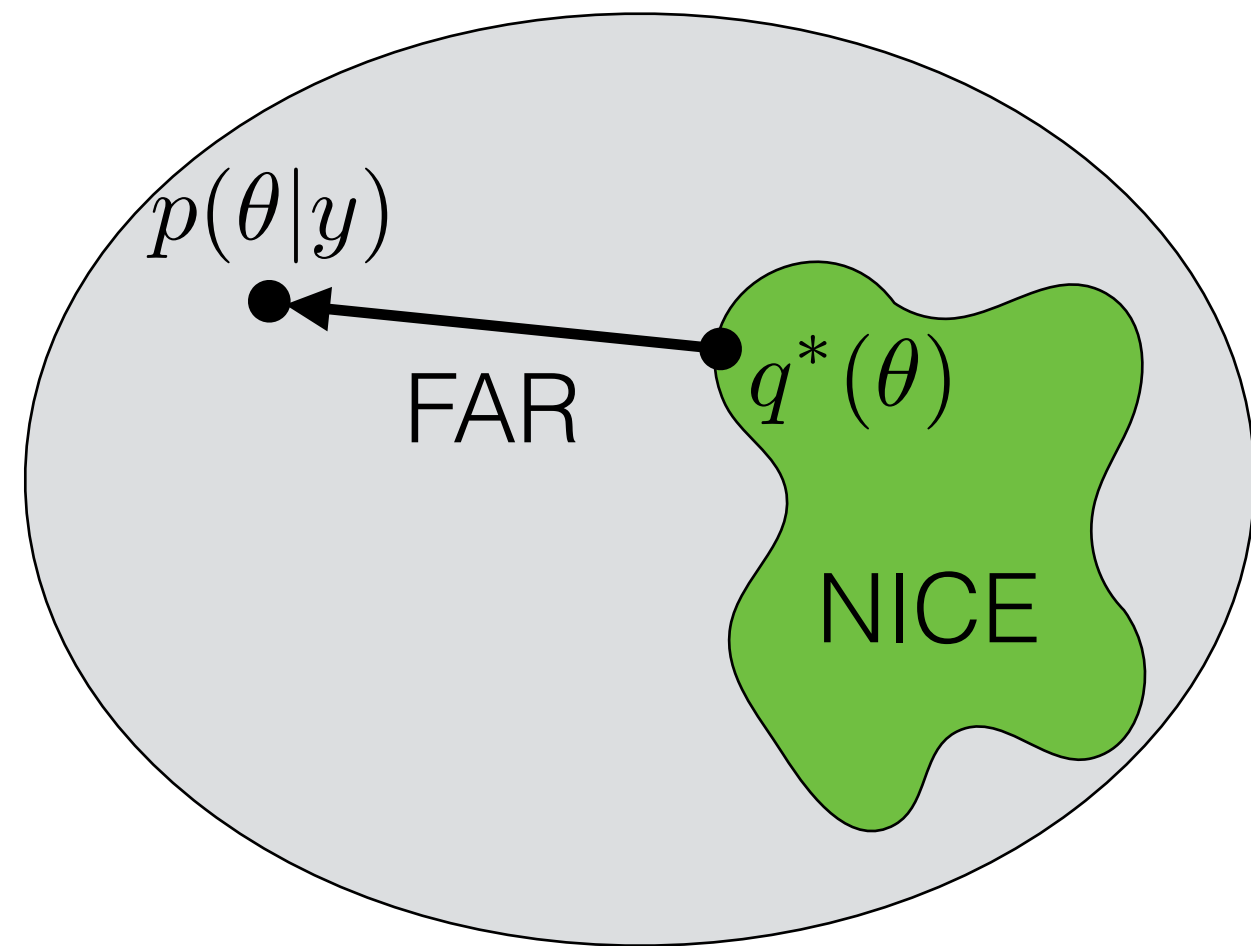
Choose “NICE” distributions



Variational Bayes

$$q^* = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

Choose “NICE” distributions



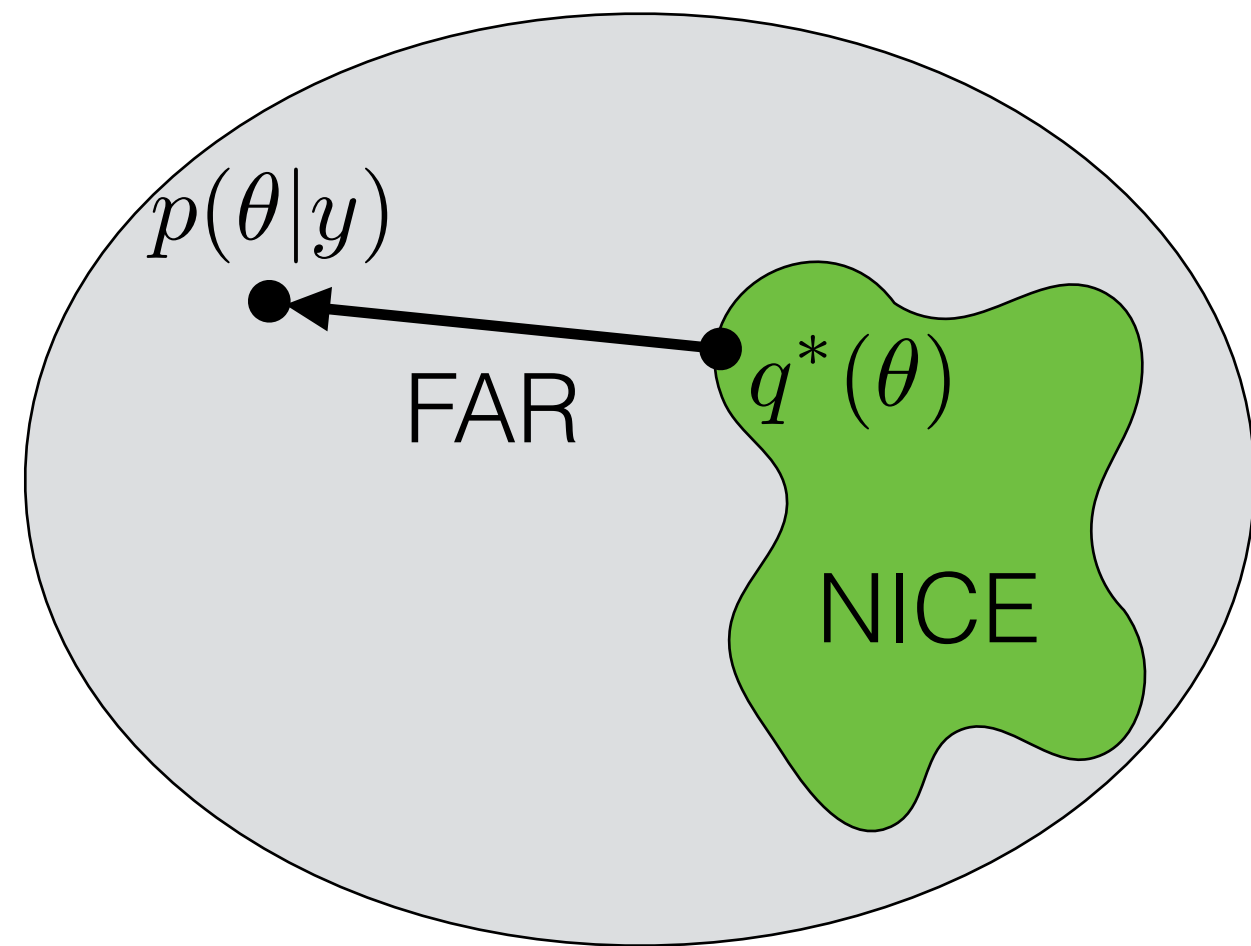
Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

Choose “NICE” distributions

- Mean-field variational Bayes (MFVB)

$$Q_{MFVB} := \left\{ q : q(\theta) = \prod_{j=1}^J q_j(\theta_j) \right\}$$



Variational Bayes

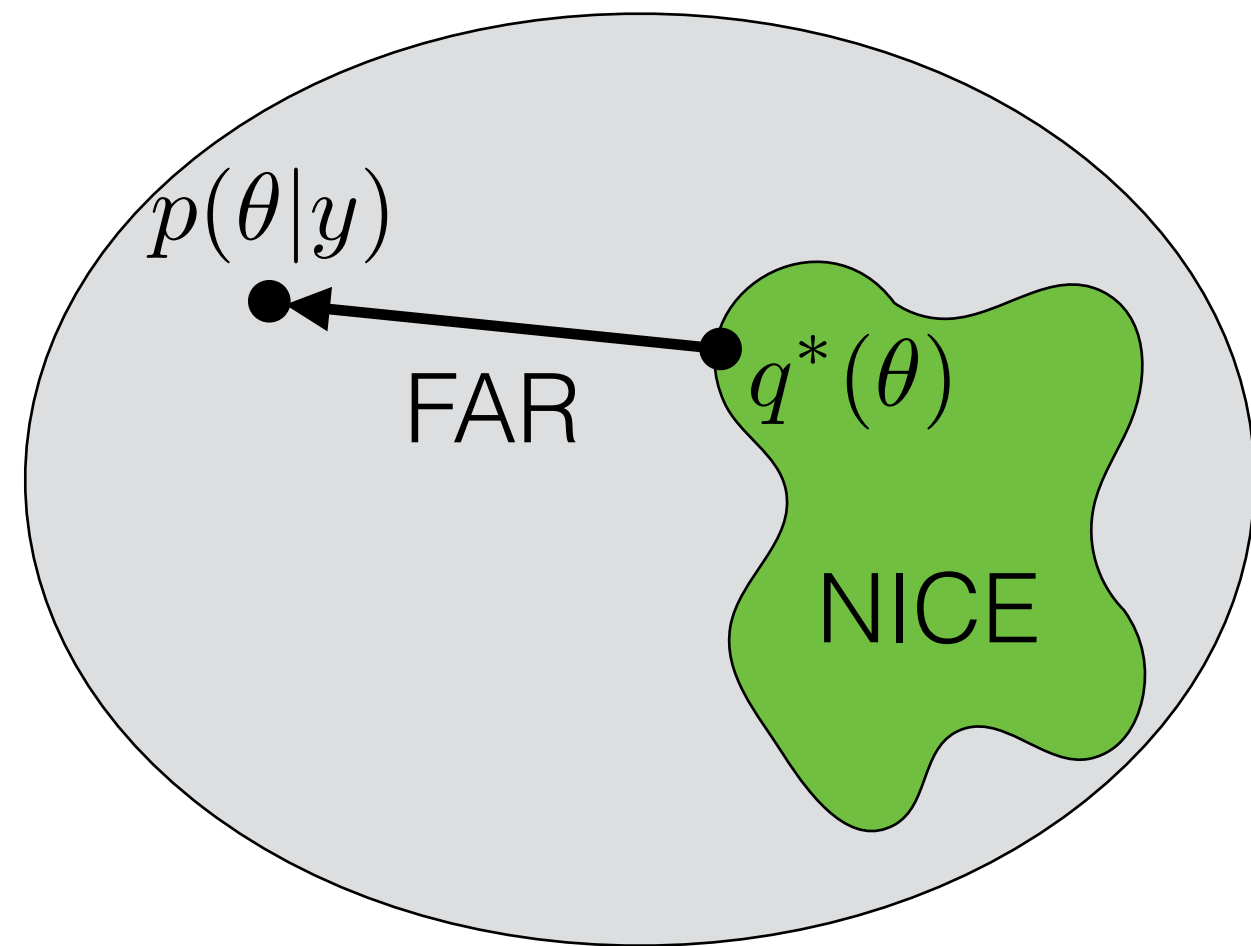
$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

Choose “NICE” distributions

- Mean-field variational Bayes (MFVB)

$$Q_{MFVB} := \left\{ q : q(\theta) = \prod_{j=1}^J q_j(\theta_j) \right\}$$

- Often also exponential family



Variational Bayes

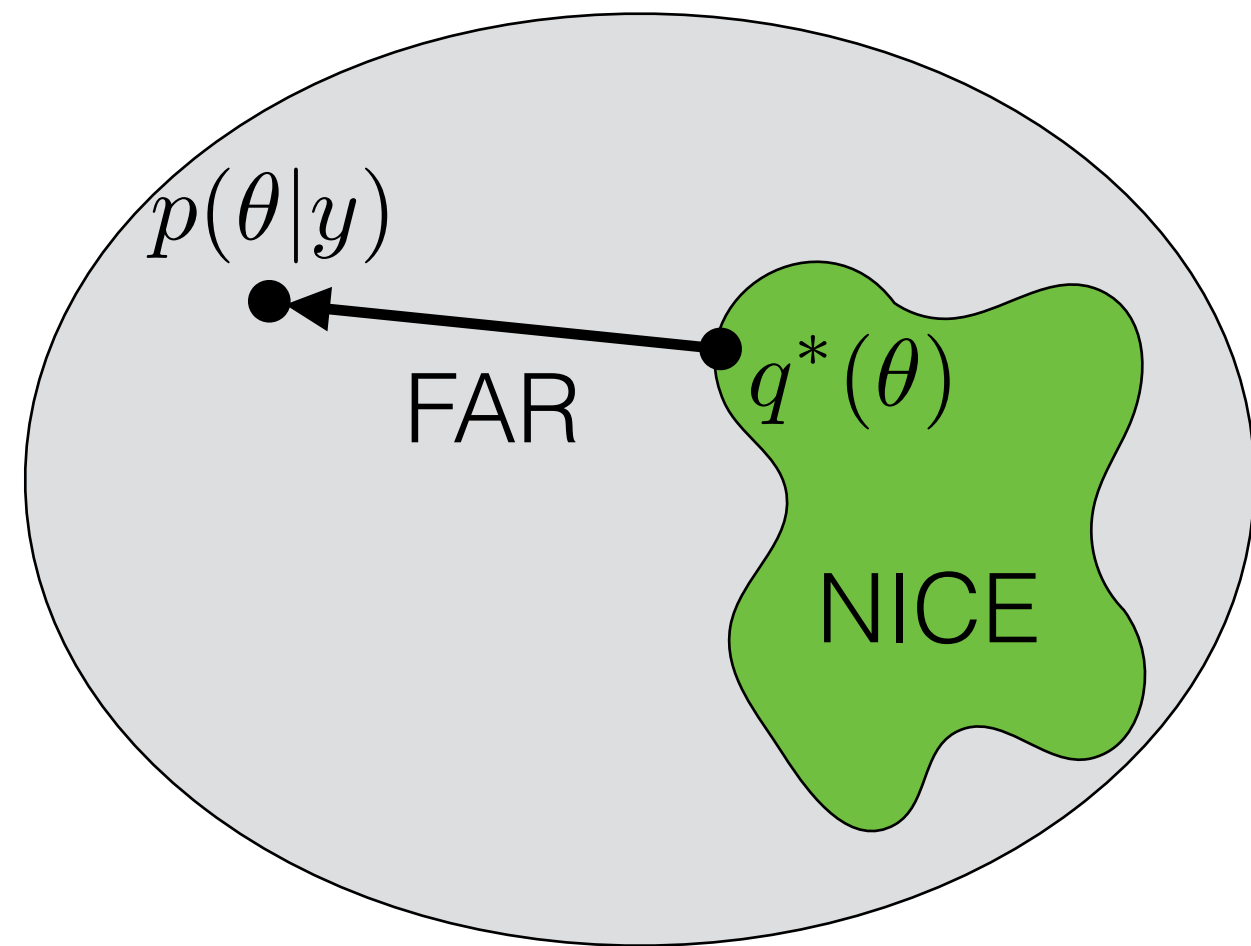
$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

Choose “NICE” distributions

- Mean-field variational Bayes (MFVB)

$$Q_{MFVB} := \left\{ q : q(\theta) = \prod_{j=1}^J q_j(\theta_j) \right\}$$

- Often also exponential family
- *Not* a modeling assumption



Variational Bayes

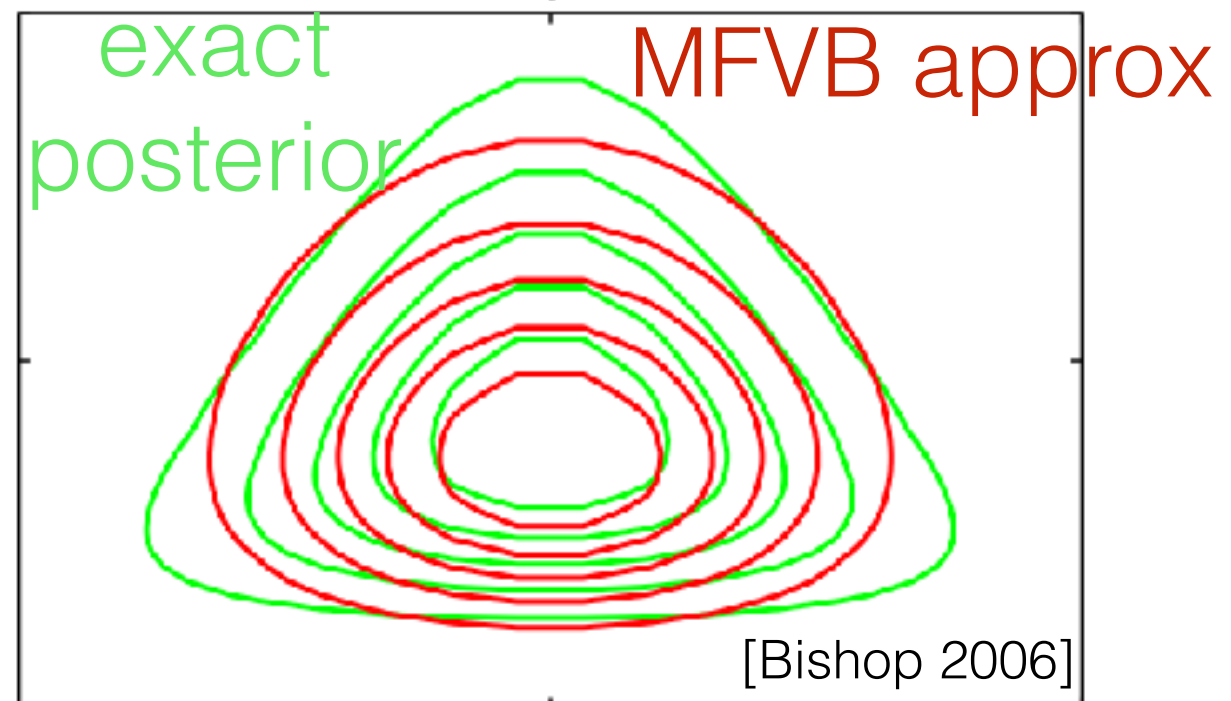
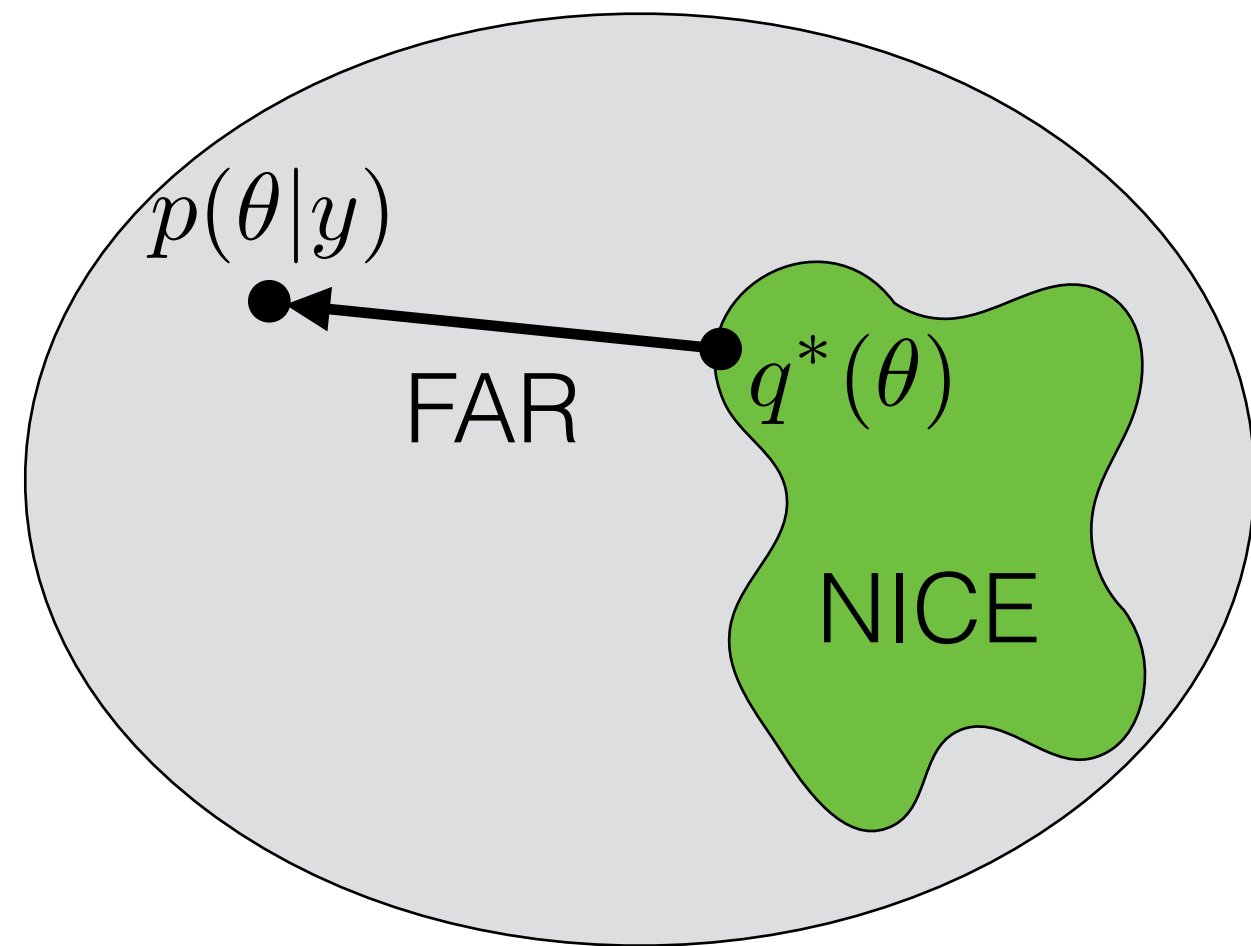
$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

Choose “NICE” distributions

- Mean-field variational Bayes (MFVB)

$$Q_{MFVB} := \left\{ q : q(\theta) = \prod_{j=1}^J q_j(\theta_j) \right\}$$

- Often also exponential family
- *Not* a modeling assumption



Variational Bayes

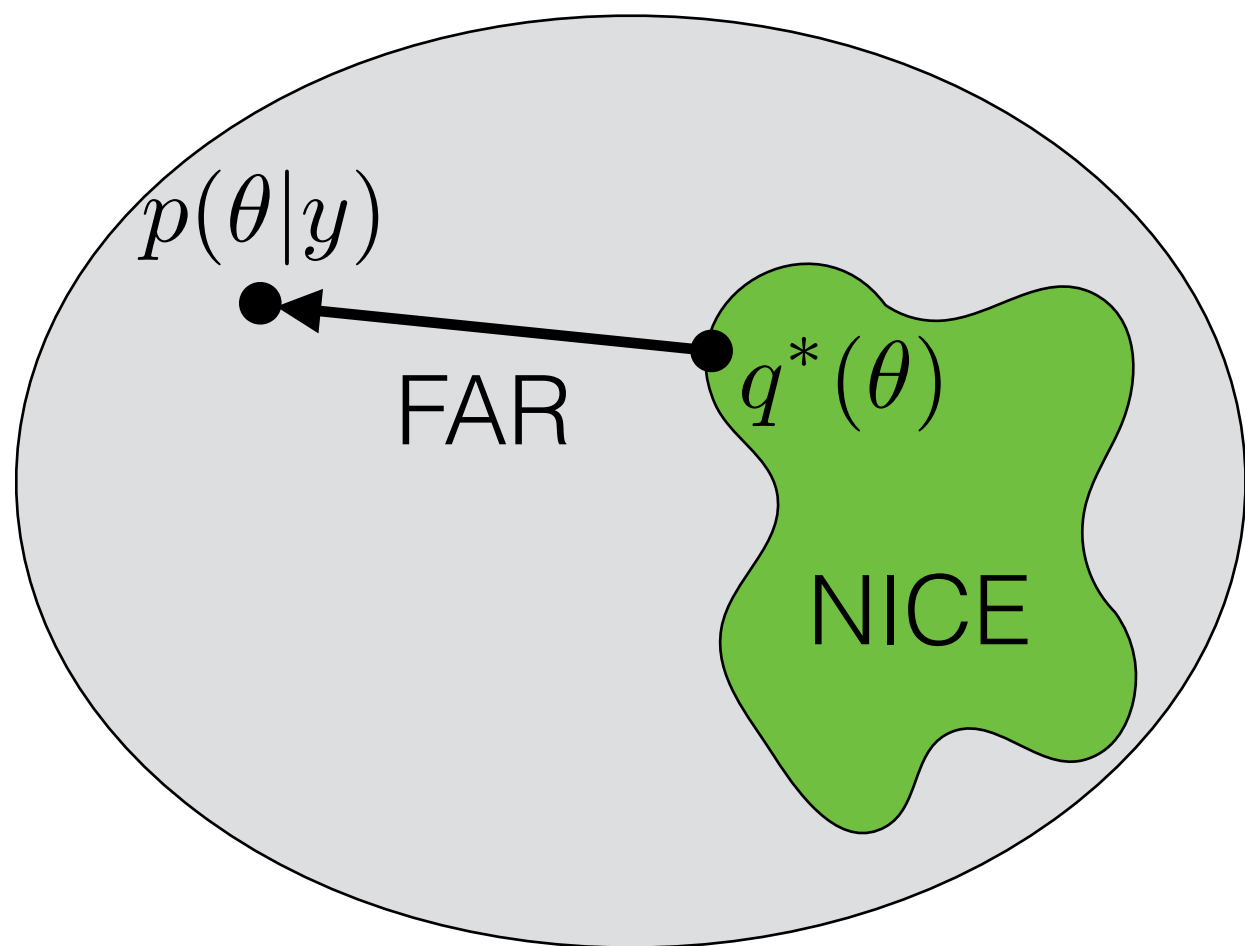
$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

Choose “NICE” distributions

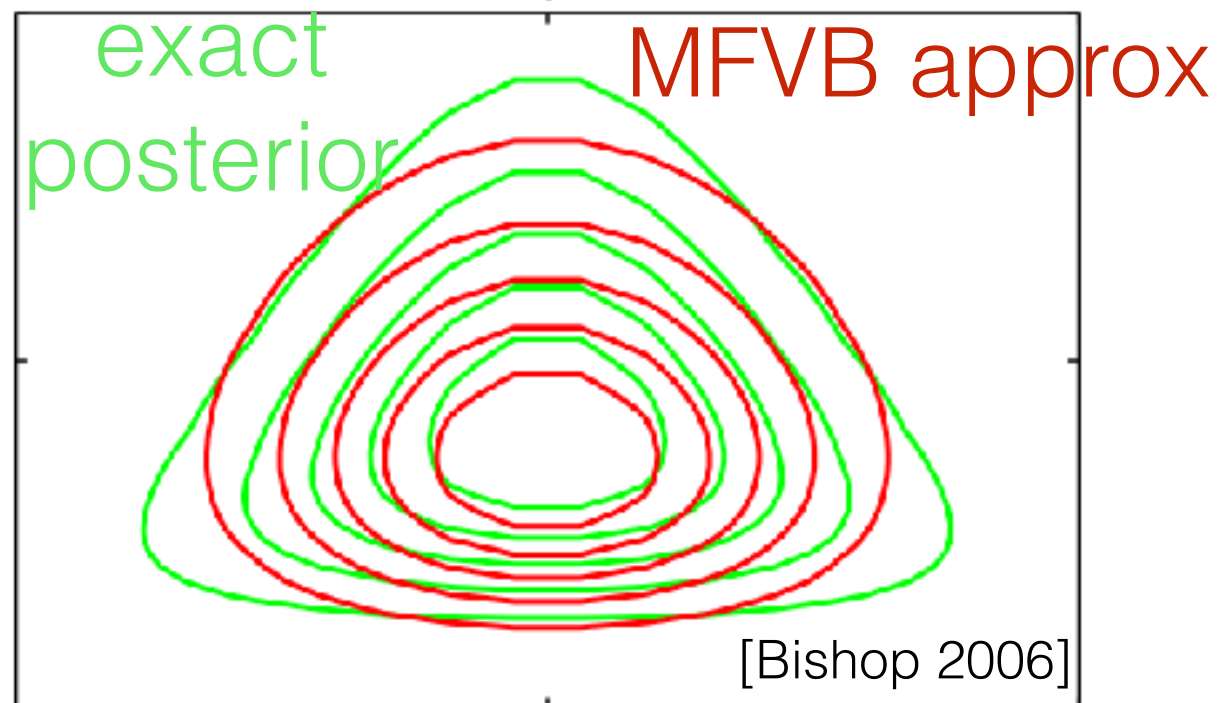
- Mean-field variational Bayes (MFVB)

$$Q_{MFVB} := \left\{ q : q(\theta) = \prod_{j=1}^J q_j(\theta_j) \right\}$$

- Often also exponential family
- *Not* a modeling assumption



Now we have an optimization problem; how to solve it?



Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} \operatorname{KL}(q(\cdot) || p(\cdot|y))$$

Choose “NICE” distributions

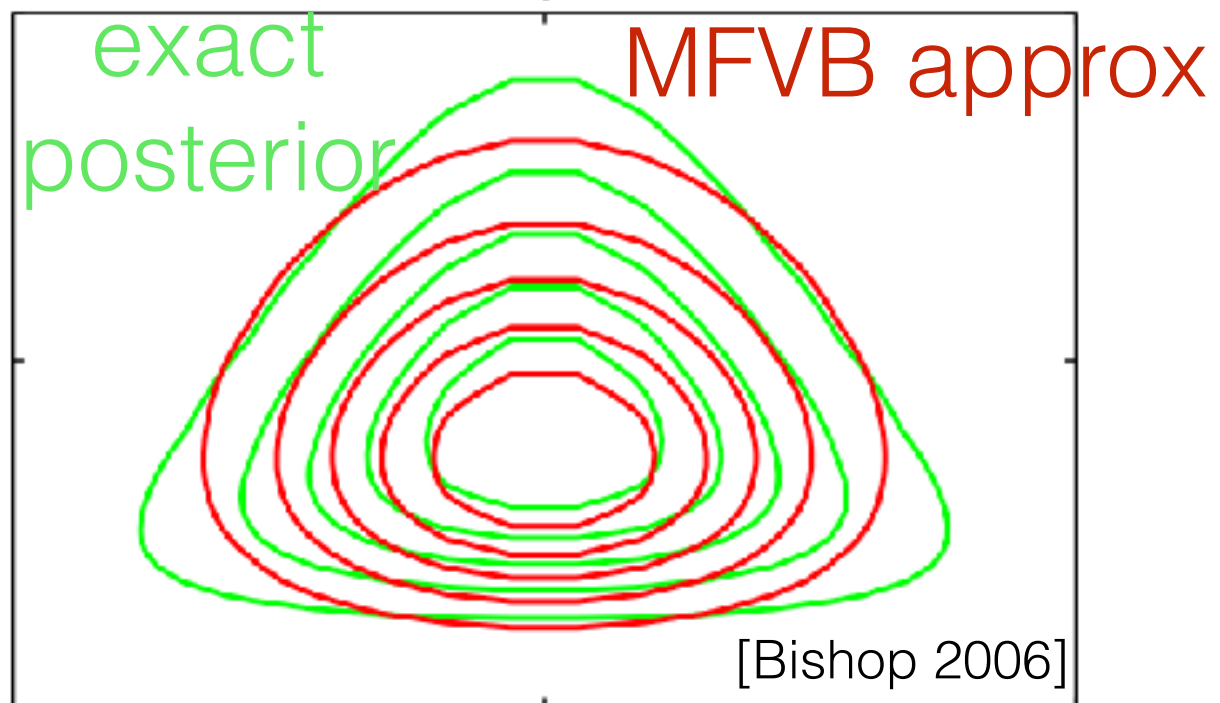
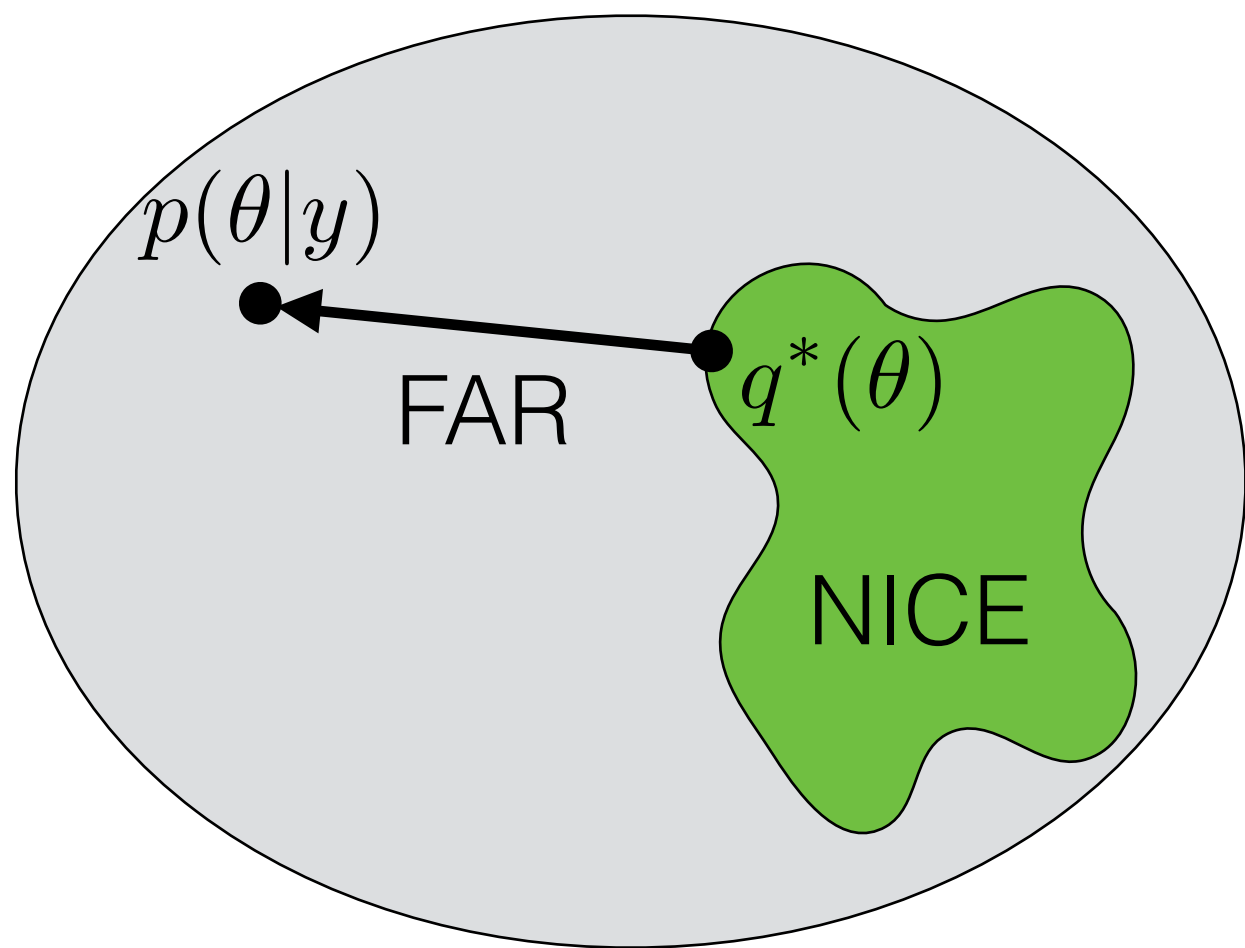
- Mean-field variational Bayes (MFVB)

$$Q_{MFVB} := \left\{ q : q(\theta) = \prod_{j=1}^J q_j(\theta_j) \right\}$$

- Often also exponential family
- *Not* a modeling assumption

Now we have an optimization problem; how to solve it?

- *One* option: Coordinate descent in q_1, \dots, q_J



Approximate Bayesian inference

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

- Coordinate descent

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

- Coordinate descent
- Stochastic variational inference (SVI) [Hoffman et al 2013]

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

- Coordinate descent
- Stochastic variational inference (SVI) [Hoffman et al 2013]
- Automatic differentiation variational inference (ADVI) [Kucukelbir et al 2015, 2017]

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

- Coordinate descent
- Stochastic variational inference (SVI) [Hoffman et al 2013]
- Automatic differentiation variational inference (ADVI) [Kucukelbir et al 2015, 2017]

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

- Coordinate descent
- Stochastic variational inference (SVI) [Hoffman et al 2013]
- Automatic differentiation variational inference (ADVI) [Kucukelbir et al 2015, 2017]

Approximate Bayesian inference

Use q^* to approximate $p(\cdot|y)$

Optimization

$$q^* = \operatorname{argmin}_{q \in Q} f(q(\cdot), p(\cdot|y))$$

Variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q} KL(q(\cdot) || p(\cdot|y))$$

Mean-field variational Bayes

$$q^* = \operatorname{argmin}_{q \in Q_{\text{MFVB}}} KL(q(\cdot) || p(\cdot|y))$$

- Coordinate descent
- Stochastic variational inference (SVI) [Hoffman et al 2013]
- Automatic differentiation variational inference (ADVI) [Kucukelbir et al 2015, 2017]

Roadmap

- Bayes & Approximate Bayes review
- What is:
 - Variational Bayes (VB)
 - Mean-field variational Bayes (MFVB)
- Why use VB?
- When can we trust VB?
- Where do we go from here?

Roadmap

- Bayes & Approximate Bayes review
- What is:
 - Variational Bayes (VB)
 - Mean-field variational Bayes (MFVB)
- Why use VB?
- When can we trust VB?
- Where do we go from here?

What to read next

Textbooks and Reviews

- Bishop. *Pattern Recognition and Machine Learning*, Ch 10. 2006.
- Blei, Kucukelbir, McAuliffe. Variational inference: A review for statisticians, *JASA* 2016.
- MacKay. *Information Theory, Inference, and Learning Algorithms*, Ch 33. 2003.
- Murphy. *Machine Learning: A Probabilistic Perspective*, Ch 21. 2012.
- Ormerod, Wand. Explaining variational approximations. *Amer Stat* 2010.
- Turner, Sahani. Two problems with variational expectation maximisation for time-series models. In *Bayesian Time Series Models*, 2011.
- Wainwright, Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 2008.

Our Experiments

- RJ Giordano, T Broderick, and MI Jordan. Linear response methods for accurate covariance estimates from mean field variational Bayes. *NeurIPS* 2015.
- RJ Giordano, T Broderick, R Meager, JH Huggins, and MI Jordan. Fast robustness quantification with variational Bayes. *ICML Data4Good Workshop* 2016.
- RJ Giordano, T Broderick, and MI Jordan. Covariances, robustness, and variational Bayes. *JMLR* 2018.
- J Huggins, M Kasprzak, T Campbell, T Broderick. Validated Variational Inference via Practical Posterior Error Bounds. ArXiv: 1910.04102. *AISTATS* 2020, to appear.
- T Campbell and T Broderick. Automated scalable Bayesian inference via Hilbert coresets. *JMLR* 2019.
- T Campbell and T Broderick. Bayesian Coreset Construction via Greedy Iterative Geodesic Ascent. *ICML* 2018.

References

Full references at end of final slides