

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Expectation Propagation in the large data limit

Guillaume Dehaene, Simon Barthelmé

February 29, 2016

Bayesian inference is hard

- Bayesian inference is a powerful statistical methodology
- But for most generative models, it's also computationally intractable

$$p(x|O_1 \dots O_n) = p(x) \prod_i p(O_i|x)$$

Bayesian inference is hard

- Bayesian inference is a powerful statistical methodology
- But for most generative models, it's also computationally intractable

$$p(x|O_1 \dots O_n) = p(x) \prod_i p(O_i|x)$$

- What can we do ?
 - Point estimates ! (maximum likelihood, MAP)

Bayesian inference is hard

- Bayesian inference is a powerful statistical methodology
- But for most generative models, it's also computationally intractable

$$p(x|O_1 \dots O_n) = p(x) \prod_i p(O_i|x)$$

- What can we do ?
 - Point estimates ! (maximum likelihood, MAP)
 - Sampling methods ! Generate $X \sim p(x|O_1 \dots O_n)$

Bayesian inference is hard

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Bayesian inference is a powerful statistical methodology
- But for most generative models, it's also computationally intractable

$$p(x|O_1 \dots O_n) = p(x) \prod_i p(O_i|x)$$

- What can we do ?
 - Point estimates ! (maximum likelihood, MAP)
 - Sampling methods ! Generate $X \sim p(x|O_1 \dots O_n)$
 - Approximate inference ! Find $q \approx p$

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Expectation Propagation

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm



- It's used to match players in skill level in Halo (Microsoft True Skill, Xbox)

EP is powerful ...

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- EP has great potential:
 - It's powerful (Kuss et al, 05; Nickish et al, 08):
 - Empirically, it gives high-quality approximations at minimal cost
 - It's universal:
 - it can be applied to any $p(x)$ with a simple factor structure
 - Can perform the computation in parallel

but poorly understood !!

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- But EP is also very poorly known !!
- Open questions:
 - How good are the approximations ?
 - Does it always terminate ?

but poorly understood !!

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- But EP is also very poorly known !!
- Open questions:
 - How good are the approximations ?
 - Does it always terminate ?
- We've been able to tackle those questions in the large-data limit:
 - We prove it gives good approximations
 - We prove that it has a simple limit behavior

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

1 Background

How does EP work ?

The large-data limit

2 Why does EP give accurate approximations

3 The EP iteration behaves like Newton's algorithm

Notations

We will approximate a 1D probability distribution $p(x)$ that has a simple factor structure

$$p(x|O_1 \dots O_n) = p(x) \prod_{i=1}^n p(O_i|x)$$

$$p(x) = \prod_{i=1}^n f_i(x)$$

Notations

We will approximate a 1D probability distribution $p(x)$ **that has a simple factor structure**

$$p(x|O_1 \dots O_n) = p(x) \prod_{i=1}^n p(O_i|x)$$

$$p(x) = \prod_{i=1}^n f_i(x)$$

We will approximate p with a Gaussian g that also factorizes:

$$g(x) = \prod_{i=1}^n g_i(x) \approx p(x)$$
$$\forall i \quad g_i(x) \approx f_i(x)$$

Notations

We will approximate a 1D probability distribution $p(x)$ **that has a simple factor structure**

$$p(x|O_1 \dots O_n) = p(x) \prod_{i=1}^n p(O_i|x)$$

$$p(x) = \prod_{i=1}^n f_i(x)$$

We will approximate p with a Gaussian g that also factorizes:

$$g(x) = \prod_{i=1}^n g_i(x) \approx p(x)$$
$$\forall i \quad g_i(x) \approx f_i(x)$$

We will often work with negative logs:

$$\psi(x) = -\log[p(x)]$$

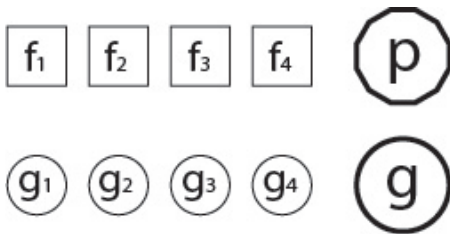
$$\phi_i(x) = -\log[f_i(x)]$$

The EP loop

- EP aims to find a factorized approximation of $p(x)$:

$$g(x) = \prod_{i=1}^n g_i(x) \approx \prod_{i=1}^n f_i(x) = p(x)$$

- EP proceeds **iteratively**
- The basic idea: How do we improve a current approximation $[g_i^t(x)]$??

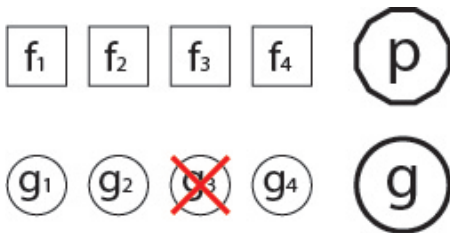


The EP loop

- EP aims to find a factorized approximation of $p(x)$:

$$g(x) = \prod_{i=1}^n g_i(x) \approx \prod_{i=1}^n f_i(x) = p(x)$$

- EP proceeds **iteratively**
- The basic idea: How do we improve a current approximation $[g_i^t(x)]$??

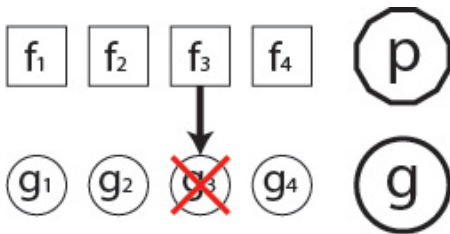


The EP loop

- EP aims to find a factorized approximation of $p(x)$:

$$g(x) = \prod_{i=1}^n g_i(x) \approx \prod_{i=1}^n f_i(x) = p(x)$$

- EP proceeds **iteratively**
- The basic idea: How do we improve a current approximation $[g_i^t(x)]$??

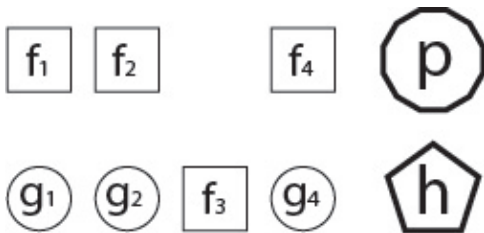


The EP loop

- EP aims to find a factorized approximation of $p(x)$:

$$g(x) = \prod_{i=1}^n g_i(x) \approx \prod_{i=1}^n f_i(x) = p(x)$$

- EP proceeds **iteratively**
- The basic idea: How do we improve a current approximation $[g_i^t(x)]$??

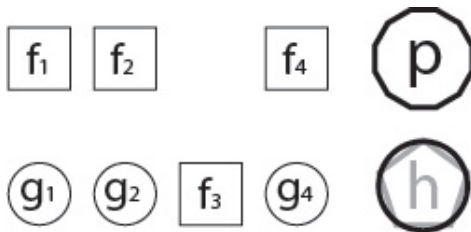


The EP loop

- EP aims to find a factorized approximation of $p(x)$:

$$g(x) = \prod_{i=1}^n g_i(x) \approx \prod_{i=1}^n f_i(x) = p(x)$$

- EP proceeds **iteratively**
- The basic idea: How do we improve a current approximation $[g_i^t(x)]$??

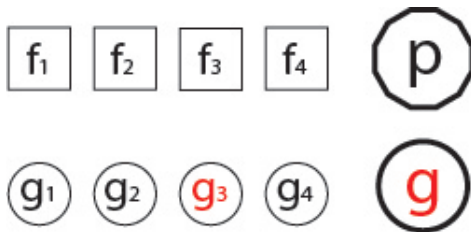


The EP loop

- EP aims to find a factorized approximation of $p(x)$:

$$g(x) = \prod_{i=1}^n g_i(x) \approx \prod_{i=1}^n f_i(x) = p(x)$$

- EP proceeds **iteratively**
- The basic idea: How do we improve a current approximation $[g_i^t(x)]$??



The EP loop

- Select i for updating
- Compute

$$h_i(x) = f_i(x) \prod_{j \neq i} g_j(x)$$

- Compute a Gaussian approximation:

$$g^{t+1}(x) \approx h_i(x)$$

- update the approximation of f_i :

$$g_i^{t+1} = \frac{g^{t+1}}{\prod_{j \neq i} g_j(x)}$$

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

The EP loop

- Select i for updating
- Compute

$$h_i(x) = f_i(x) \prod_{j \neq i} g_j(x)$$

- Compute a Gaussian approximation:

$$g^{t+1}(x) \approx h_i(x)$$

- update the approximation of f_i :

$$g_i^{t+1} = \frac{g^{t+1}}{\prod_{j \neq i} g_j(x)}$$

- Terminology:
 - $g_{-i} = \prod g_j$ is the **cavity distribution**
 - h_i is the **hybrid**

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Approximating the hybrid

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- How do we compute $g^{t+1} \approx h_i$?

Approximating the hybrid

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- How do we compute $g^{t+1} \approx h_i$?
- Minimize the **Kullback-Leibler** divergence

$$g^{t+1} = \operatorname{argmin}_g KL(h_i, g)$$

- Gives a good approximation
- Is simple to compute

Minimizing KL

- Inside exponential families, minimizing KL is easy

Background

**How does EP
work ?**

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Minimizing KL

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Inside exponential families, minimizing KL is easy
- Gaussians are an exponential family:

$$g(x|r, \beta) \propto \exp\left(rx - \beta \frac{x^2}{2}\right)$$

Minimizing KL

- Inside exponential families, minimizing KL is easy
- Gaussians are an exponential family:

$$g(x|r, \beta) \propto \exp\left(rx - \beta \frac{x^2}{2}\right)$$

- Relation between r, β and the moments:

$$\begin{aligned}\mu &= \frac{r}{\beta} \\ \text{var} &= \beta^{-1}\end{aligned}$$

Minimizing KL

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

To find $\operatorname{argmin}_g KL(h_i, g)$

- Compute the mean and variance of h_i
- Compute the Gaussian with that mean and variance:

$$r = \operatorname{var}^{-1} \mu$$

$$\beta = \operatorname{var}^{-1}$$

Working in natural parameters

- Working in the space of Gaussians: $g_i \in \mathcal{G}$: impossible to visualize

Background

**How does EP
work ?**

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Working in natural parameters

- Working in the space of Gaussians: $g_i \in \mathcal{G}$: impossible to visualize
- Working in moments:

$$\begin{aligned}\mu_i &= E_{g_i}(x) \\ v_i &= \text{var}_{g_i}(x)\end{aligned}$$

Better, but hard to multiply and divide Gaussians

Working in natural parameters

- Working in the space of Gaussians: $g_i \in \mathcal{G}$: impossible to visualize
- Working in moments:

$$\begin{aligned}\mu_i &= E_{g_i}(x) \\ \nu_i &= \text{var}_{g_i}(x)\end{aligned}$$

Better, but hard to multiply and divide Gaussians

- Working in natural parameters:

$$g_i(x) \propto \exp\left(r_i x - \beta_i \frac{x^2}{2}\right)$$

Multiplication and division of Gaussians = sums and differences of natural parameters !

EP in natural parameters

Sequential algorithm, operating on $(2n)$ dimensional space
 $[r_i, \beta_i]$

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- For i in $1 \dots n$

- ① Compute “cavity” parameters: $r_{-i} = \sum_{j \neq i} r_j$, $\beta_{-i} = \sum_{j \neq i} \beta_j$
- ② Compute hybrid distribution $h_i(x) = f_i(x) \mathcal{N}(x | r_{-i}, \beta_{-i})$
- ③ Compute $E_{h_i}(x)$ and $\text{var}_{h_i}(x)$
- ④ Update r_i and β_i from the moments of the hybrid

$$r_i = \frac{E_{h_i}(x)}{\text{var}_{h_i}} - r_{-i}$$
$$\beta_i = \frac{1}{\text{var}_{h_i}} - \beta_{-i}$$

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

An example !

Background

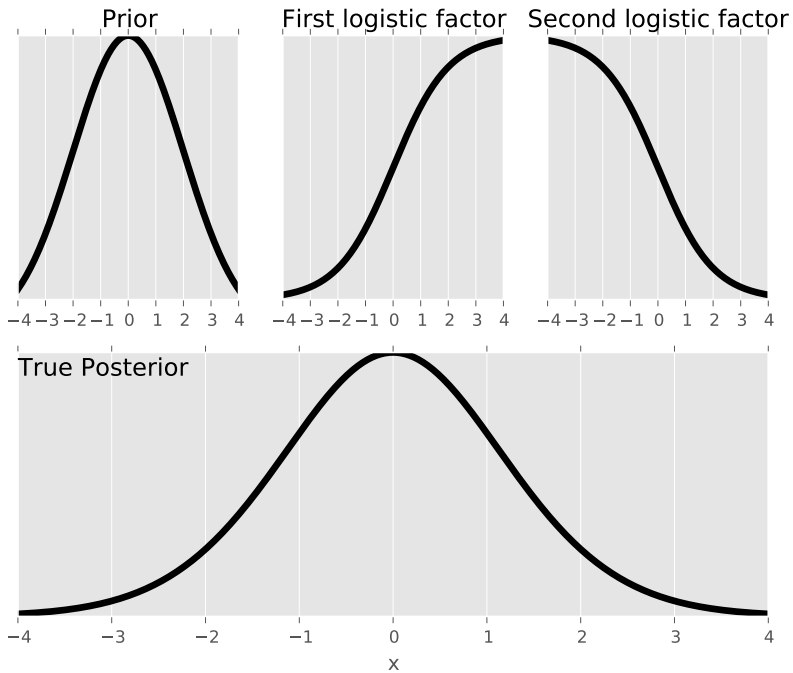
How does EP
work ?

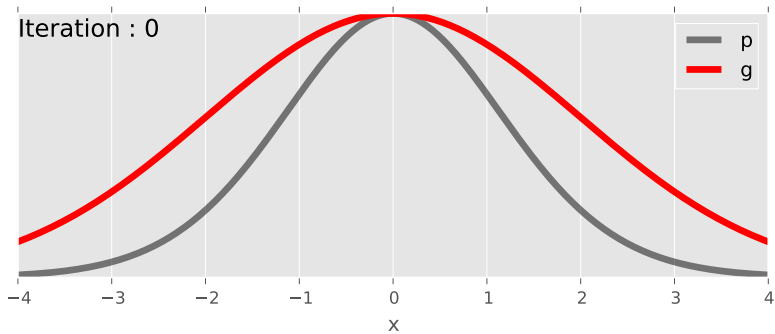
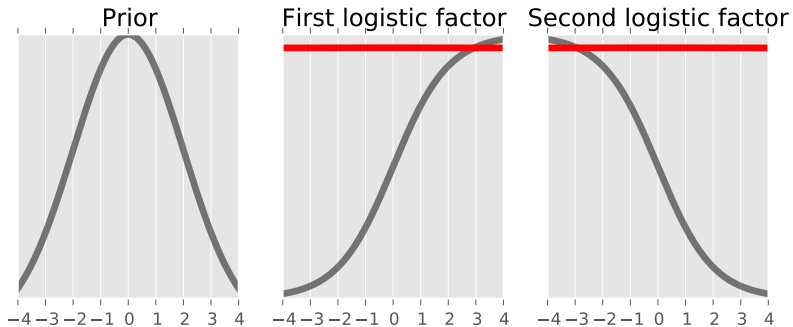
The
large-data
limit

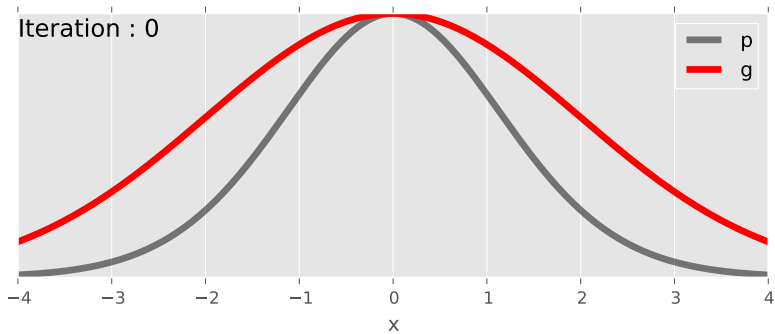
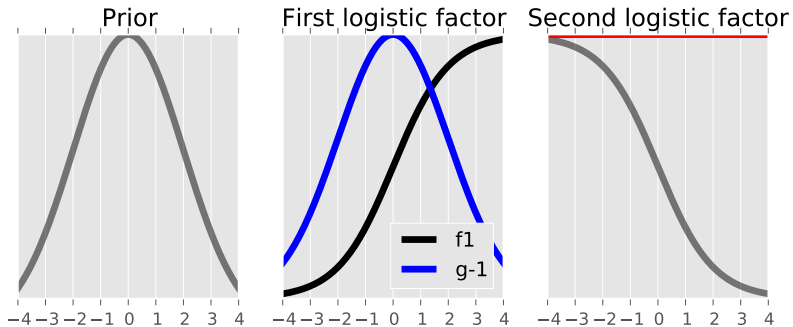
Why does
EP give
accurate ap-
proximations

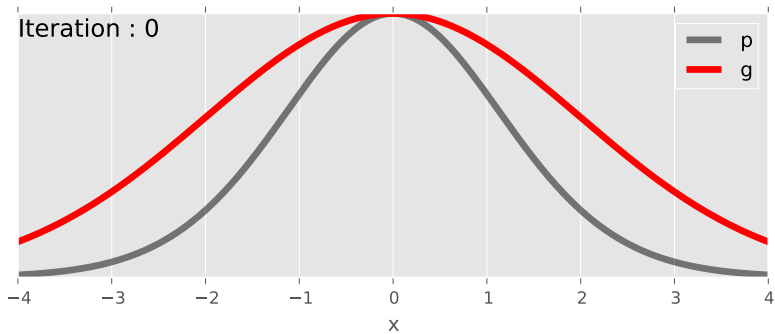
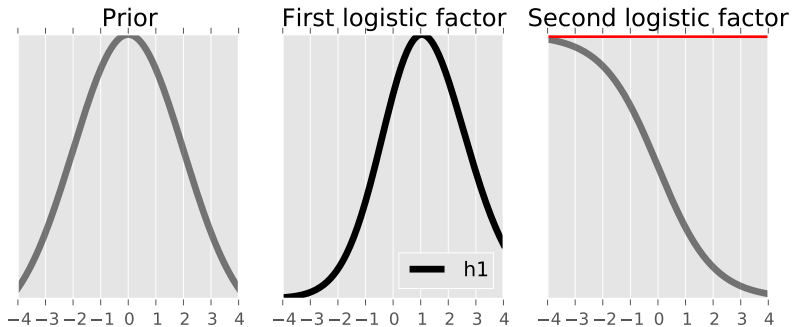
The EP
iteration
behaves like
Newton's
algorithm

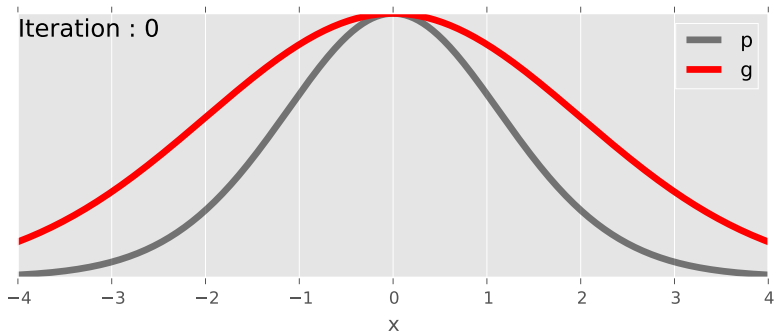
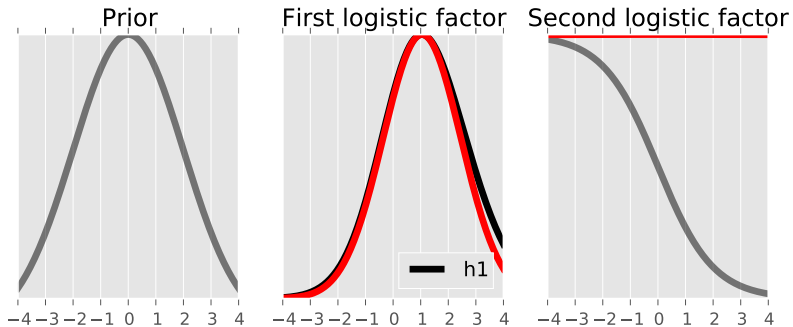
- 3 factors $f_i(x)$:
 - 2 logistic (likelihoods)
 - 1 Gaussian (prior)

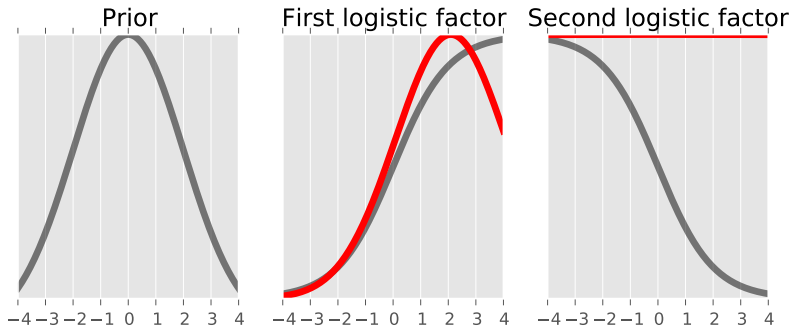


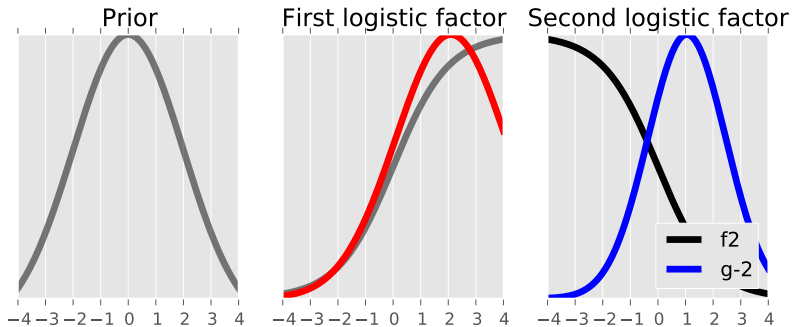


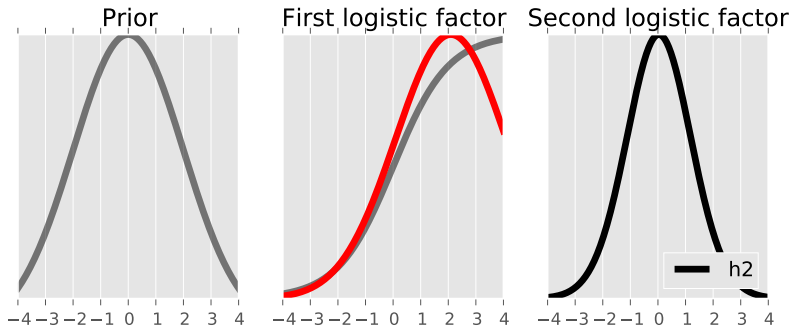


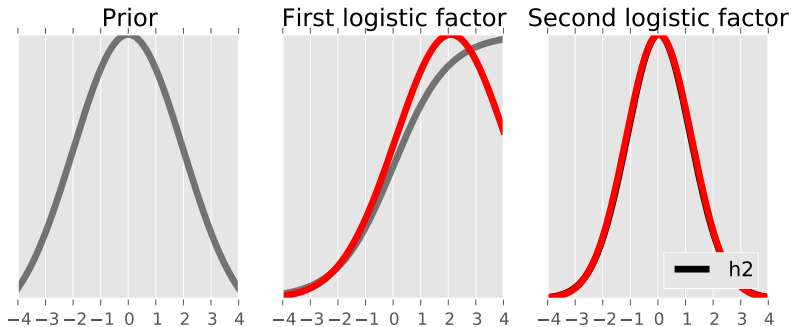


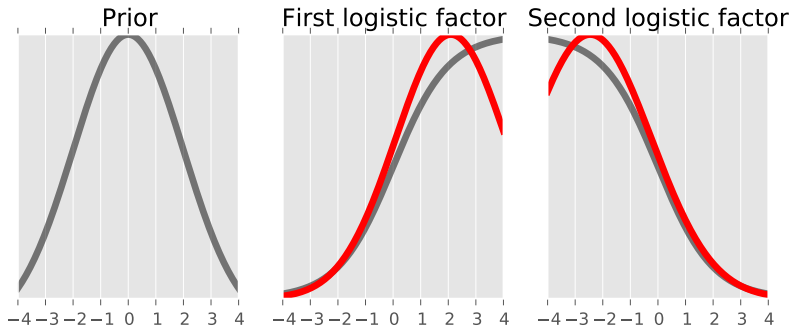


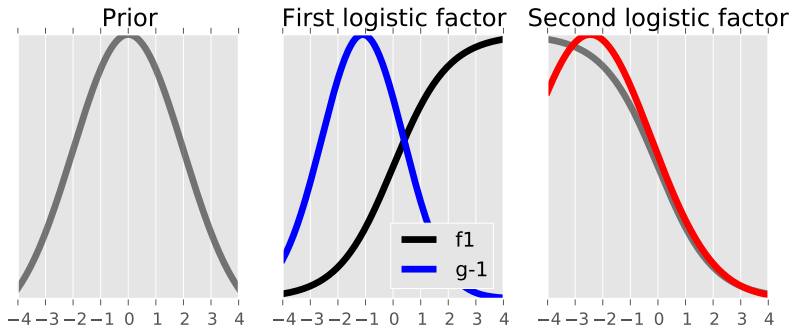


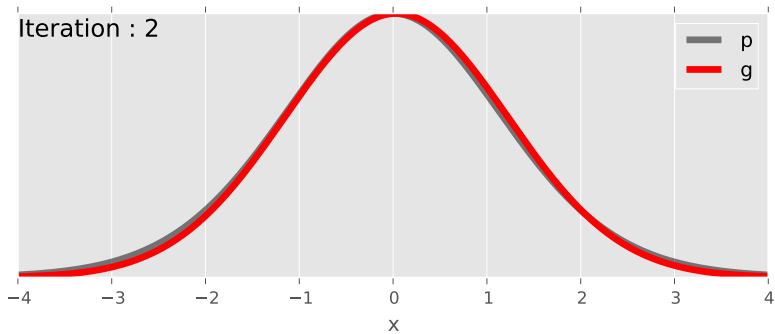
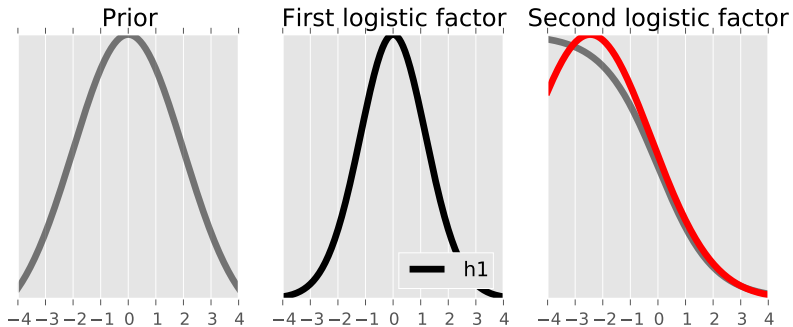


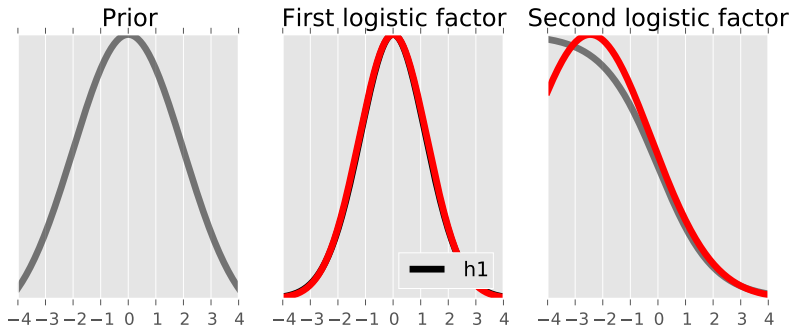


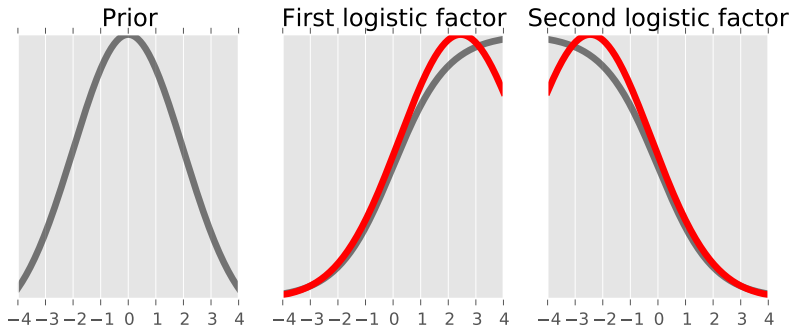


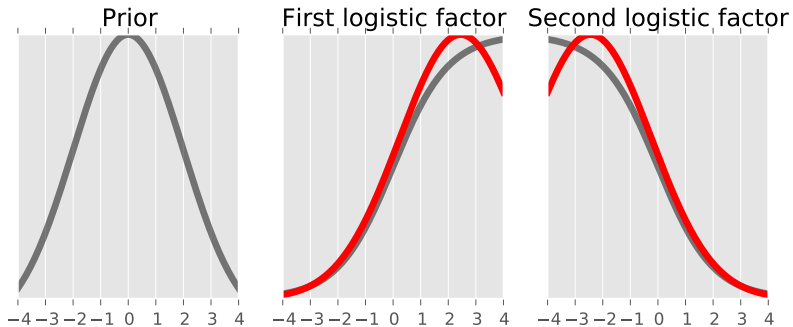


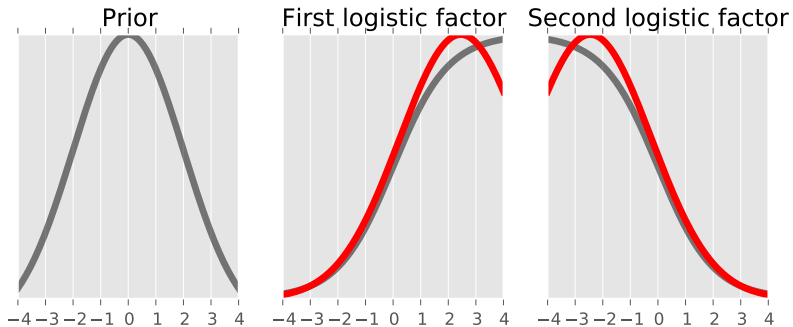


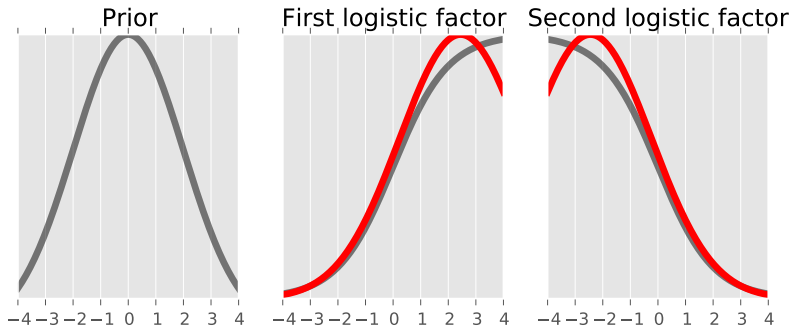


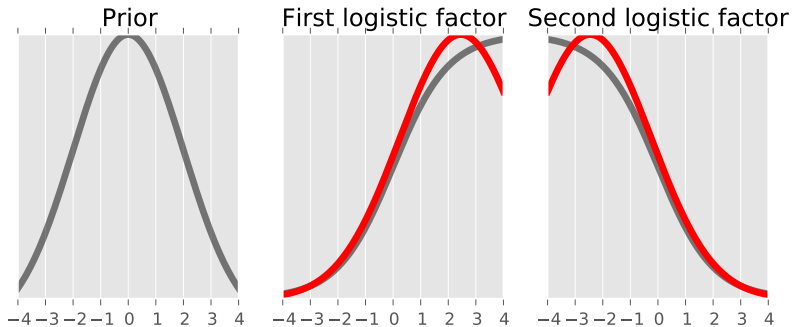












EP variants

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We presented the base algorithm which is **sequential**:
 - Pick i , then update g_i , ...

EP variants

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We presented the base algorithm which is **sequential**:
 - Pick i , then update g_i , ...
- We can also:
 - Update all approximations at once (**parallel EP**)
 - Update 10% (batch EP)
 - Update asynchronously

EP variants

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We presented the base algorithm which is **sequential**:
 - Pick i , then update g_i , ...
- We can also:
 - Update all approximations at once (**parallel EP**)
 - Update 10% (batch EP)
 - Update asynchronously
- We can also “slow-down” the updates

EP variants

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We presented the base algorithm which is **sequential**:
 - Pick i , then update g_i , ...
- We can also:
 - Update all approximations at once (**parallel EP**)
 - Update 10% (batch EP)
 - Update asynchronously
- We can also “slow-down” the updates
- All of those don't modify the fixed-points !

EP summary

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- EP approximates each factor f_i as a Gaussian g_i , and refines these approximations iteratively
- $h_i = f_i g_{-i}$ is a better approximation than g

EP summary

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- EP approximates each factor f_i as a Gaussian g_i , and refines these approximations iteratively
- $h_i = f_i g_{-i}$ is a better approximation than g
- The parameter space is the natural parameters: (r_i, β_i)
- Variants of EP modify the updating schedule, or change the updating rule

When to use EP

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- EP **apparently** works well if all h_i are almost Gaussian
- EP is dangerous to use on multimodal distributions. Like VB, EP sometimes fits a single mode of $p(x)$, missing most of the probability mass
- The EP iteration can be frustrating:
 - **slow it down** or **do it sequentially**

The large-data limit

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Large-data limit: number of observations tends to ∞
- Frequentist result: **Central Limit Theorem**: the distribution of empirical means become Gaussian

The large-data limit

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Large-data limit: number of observations tends to ∞
- Frequentist result: **Central Limit Theorem**: the distribution of empirical means become Gaussian
- Bayesian result: **Bernstein-von Mises**: posteriors converge to Gaussian distributions
- And the variance quickly goes to 0:

$$\text{var}_p(x) \propto n^{-1}$$

The large-data limit and EP

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- If approximate inference methods aren't exact in the large-data limit, they shouldn't be used

The large-data limit and EP

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- If approximate inference methods aren't exact in the large-data limit, they shouldn't be used
- The large-data limit makes theoretical analysis simple
 - The influence of a single factor f_i becomes negligible

The large-data limit and EP

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- If approximate inference methods aren't exact in the large-data limit, they shouldn't be used
- The large-data limit makes theoretical analysis simple
 - The influence of a single factor f_i becomes negligible
 - In the hybrid distribution, $h_i = f_i g_{-i}$, the cavity dominates

Contents

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

1 Background

How does EP work ?

The large-data limit

2 Why does EP give accurate approximations

3 The EP iteration behaves like Newton's algorithm

EP gives very good approximations

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We know empirically that fixed-points of EP give very good approximations of $p(x)$:
 - can we prove it ?

Assumptions

- We will constrain the factors $f_i(x)$

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Assumptions

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We will constrain the factors $f_i(x)$
- We will assume that all $f_i \propto \exp(-\phi_i)$ are **strongly log-concave**:

$$\phi_i''(x) \geq \beta_m$$

- This is an unrealistic assumption

Assumptions

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We will constrain the factors $f_i(x)$
- We will assume that all $f_i \propto \exp(-\phi_i)$ are **strongly log-concave**:

$$\phi_i''(x) \geq \beta_m$$

- This is an unrealistic assumption
- We will assume that the higher-derivatives are bounded:

$$\left| \phi_i^{(d)}(x) \right| \leq K_d$$

Assumptions

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We will constrain the factors $f_i(x)$
- We will assume that all $f_i \propto \exp(-\phi_i)$ are **strongly log-concave**:

$$\phi_i''(x) \geq \beta_m$$

- This is an unrealistic assumption
- We will assume that the higher-derivatives are bounded:

$$\left| \phi_i^{(d)}(x) \right| \leq K_d$$

- These assumptions transfer from the f_i to $p = \prod f_i$

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

The “Laplace” approximation

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We will compare EP fixed-points to the “Laplace” approximation (LA):

The “Laplace” approximation

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- We will compare EP fixed-points to the “Laplace” approximation (LA):
- Find the mode x^* of $p(x)$
- At x^* , compute $\psi''(x^*)$

$$p(x) \approx \exp \left(-\psi''(x^*) \frac{(x - x^*)^2}{2} \right)$$

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Why LA is good

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- The Bernstein-von Mises theorem justifies the LA:
 - In the large-data limit, $p_n(x) \rightarrow g_{LA}(x)$

Why LA is good

Background

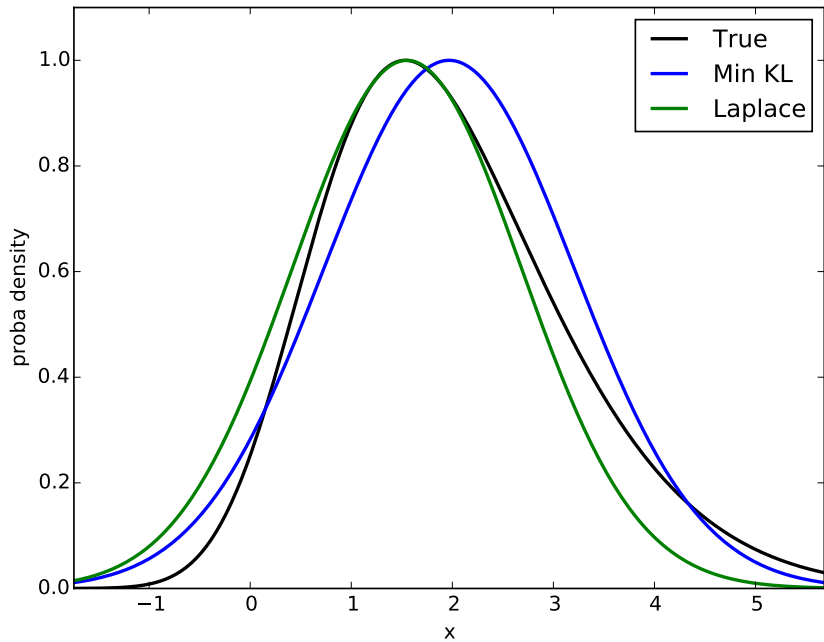
How does EP
work ?

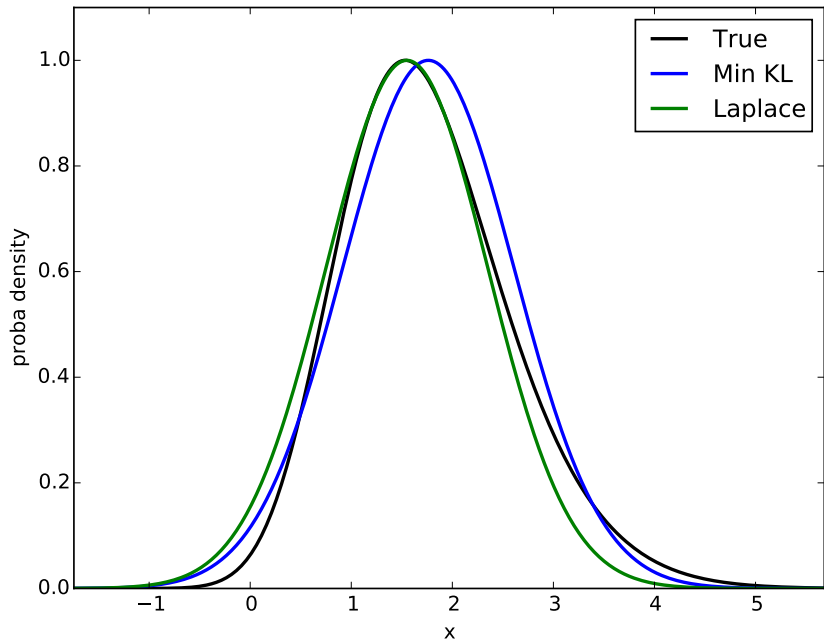
The
large-data
limit

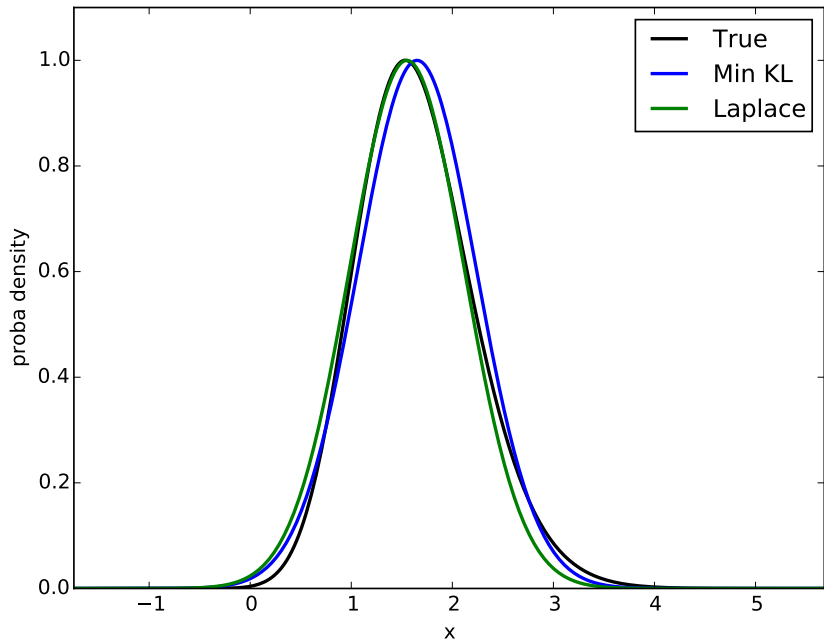
Why does
EP give
accurate ap-
proximations

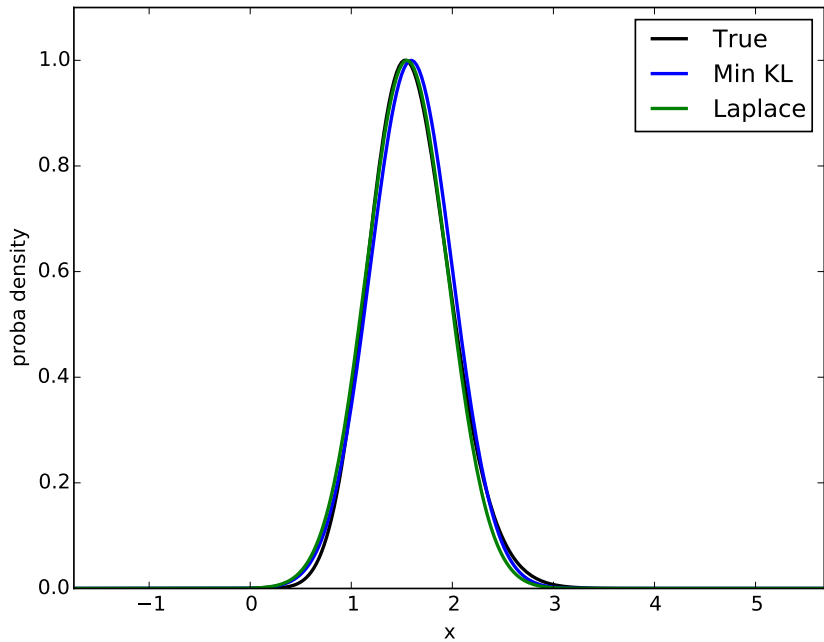
The EP
iteration
behaves like
Newton's
algorithm

- The Bernstein-von Mises theorem justifies the LA:
 - In the large-data limit, $p_n(x) \rightarrow g_{LA}(x)$
- But that doesn't mean it's perfect:
 - it looks at point estimates
 - it ignores higher derivatives









Why LA is good

- We can derive the expression of the bias:

$$x^* - \mu = -\frac{\psi^{(3)}(x^*)}{\psi''(x^*)^2} + O(n^{-2})$$

$$\psi''(x^*) - v = O(n^{-2})$$

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Why LA is good

- We can derive the expression of the bias:

$$\begin{aligned}x^{\star} - \mu &= -\frac{\psi^{(3)}(x^{\star})}{\psi''(x^{\star})^2} + O(n^{-2}) \\ \psi''(x^{\star}) - \nu &= O(n^{-2})\end{aligned}$$

- Since LA misses the mean consistently, there is **room for improvement**
- If EP is able to always correct this miss, it will improve on LA

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Why EP is better !!

Background

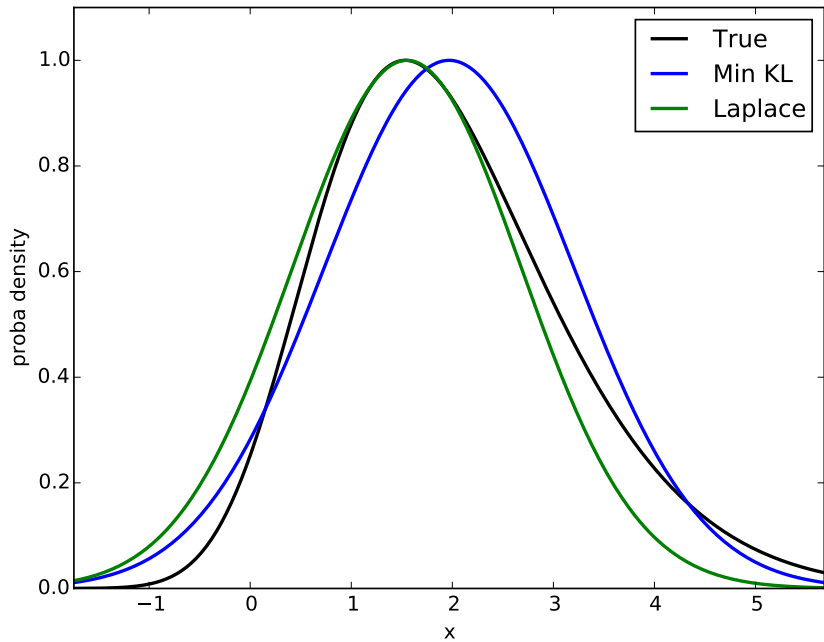
How does EP
work ?

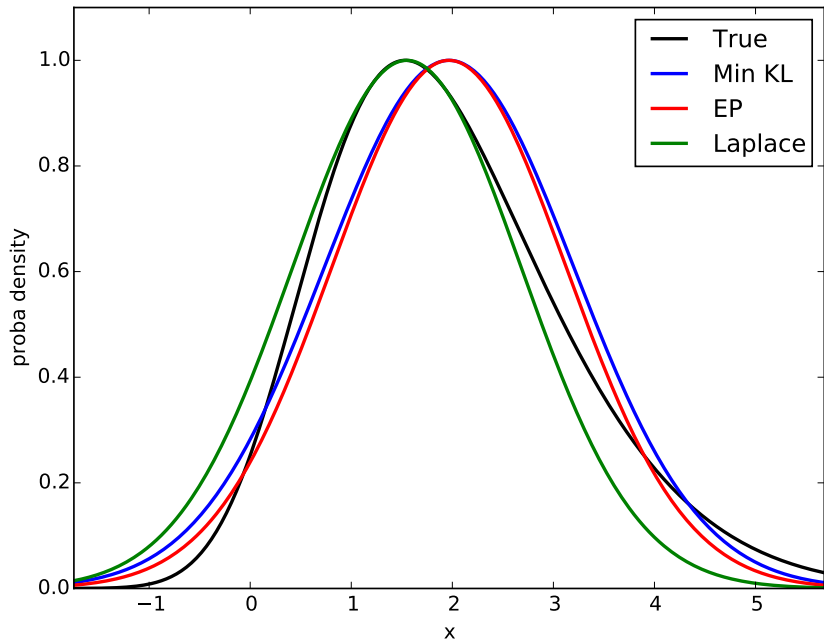
The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Consider an EP fixed-point $g(x) = \prod g_i(x) \approx p(x)$
- EP captures the $\psi^{(3)}(x^*)$ deviation





Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Why EP is better !!

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- With a similar proof as for the LA result, we prove:

$$\mu_{EP} - \mu = O(n^{-2})$$

$$v_{EP} - v = O(n^{-2})$$

Comparing LA and EP

Theorem (Quality of the LA and EP approximations)

- *LA*:

$$\begin{aligned}\mu - x^* &= O(n^{-1}) \\ \nu - \left[\psi''(x^*) \right]^{-1} &= O(n^{-2})\end{aligned}$$

- *EP*:

$$\begin{aligned}\mu - \mu_{EP} &= O(n^{-2}) \\ \nu - \nu_{EP} &= O(n^{-2})\end{aligned}$$

- *The first term of the error for the variance is slightly smaller for EP than for the LA*

The log-concavity assumption

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- The strongly log-concave sites assumption is unrealistic
- However, simple log-concavity should be enough
- proof ?

Summary

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Both EP and LA give asymptotically correct approximations of $p(x)$
- But LA fails slightly on asymmetric distributions whereas EP doesn't

Summary

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Both EP and LA give asymptotically correct approximations of $p(x)$
- But LA fails slightly on asymmetric distributions whereas EP doesn't
- Important result for credible intervals from EP approximations
- **But** problematic assumptions

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- 1 Background
 - How does EP work ?
 - The large-data limit
- 2 Why does EP give accurate approximations
- 3 The EP iteration behaves like Newton's algorithm

Understanding the EP iteration

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- The EP iteration has **one complicated step**: the site-approximation update
 - ① Compute hybrid distribution $h_i = f_i g_{-i}$
 - ② Compute $E_{h_i}(x)$ and var_{h_i}
 - ③ Compute the Gaussian with same mean and variance:
 $\mathcal{N}(x | E_{h_i}(x); var_{h_i})$
 - ④ update g_i

Understanding the EP iteration

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- The EP iteration has **one complicated step**: the site-approximation update
 - ① Compute hybrid distribution $h_i = f_i g_{-i}$
 - ② Compute $E_{h_i}(x)$ and var_{h_i}
 - ③ Compute the Gaussian with same mean and variance: $\mathcal{N}(x | E_{h_i}(x); var_{h_i})$
 - ④ update g_i
- This is the step we need to understand

Assumptions

- Much looser assumptions

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Assumptions

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Much looser assumptions
- We bound the range of the second derivatives

$$\forall i, \max \left(\phi_i'' \right) - \min \left(\phi_i'' \right) \leq B$$

- Still uniform bound on the higher derivatives

$$\left| \phi_i^{(d)}(x) \right| \leq K_d$$

Assumptions

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Much looser assumptions
- We bound the range of the second derivatives

$$\forall i, \max \left(\phi_i'' \right) - \min \left(\phi_i'' \right) \leq B$$

- Still uniform bound on the higher derivatives

$$\left| \phi_i^{(d)}(x) \right| \leq K_d$$

- Applies to any GLM, can be extended so that B and K_d depend on n

Intuitive understanding

- Rewrite KL minization:

$$\begin{aligned} g_i &= [g_{-i}]^{-1} \operatorname{argmin}_g KL(h_i, g) \\ &\approx \operatorname{argmin}_g \int h_i [\log(f_i) - \log(g_i)] \end{aligned}$$

- h_i tells us where g_i needs to fit f_i

Intuitive understanding

- Rewrite KL minimization:

$$\begin{aligned} g_i &= [g_{-i}]^{-1} \operatorname{argmin}_g KL(h_i, g) \\ &\approx \operatorname{argmin}_g \int h_i [\log(f_i) - \log(g_i)] \end{aligned}$$

- h_i tells us where g_i needs to fit f_i
- If g_{-i} has very small variance (ie: β_{-i} is big):
 - g_{-i} is almost Dirac
 - $h_i \approx g_{-i}$ and is also Dirac

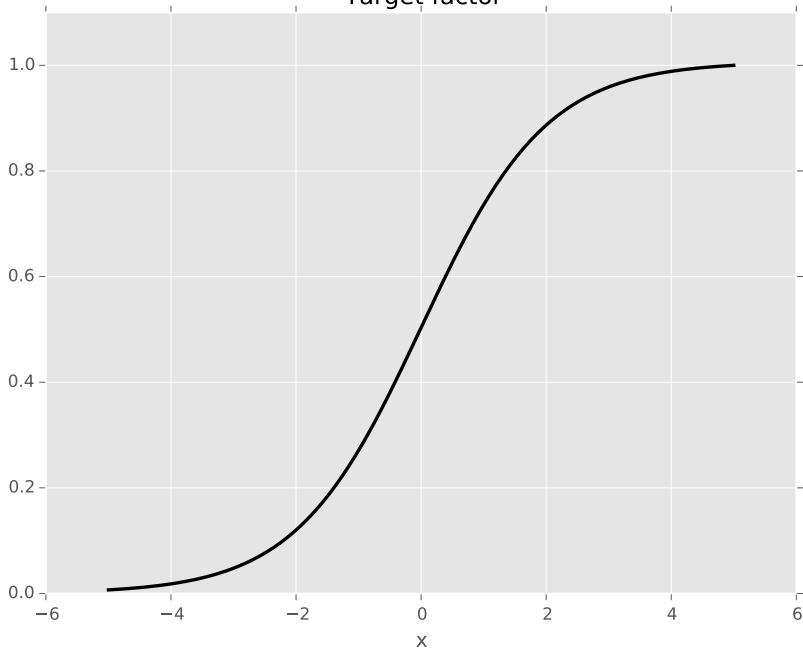
Intuitive understanding

- Rewrite KL minimization:

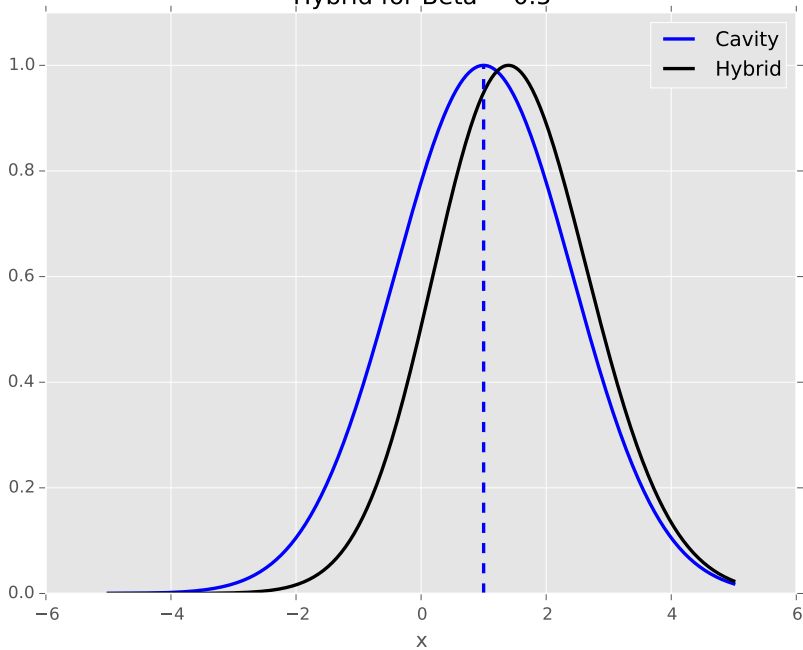
$$\begin{aligned}g_i &= [g_{-i}]^{-1} \operatorname{argmin}_g KL(h_i, g) \\ &\approx \operatorname{argmin}_g \int h_i [\log(f_i) - \log(g_i)]\end{aligned}$$

- h_i tells us where g_i needs to fit f_i
- If g_{-i} has very small variance (ie: β_{-i} is big):
 - g_{-i} is almost Dirac
 - $h_i \approx g_{-i}$ and is also Dirac
- The best approximation is the Taylor expansion of $\log(f_i)$.

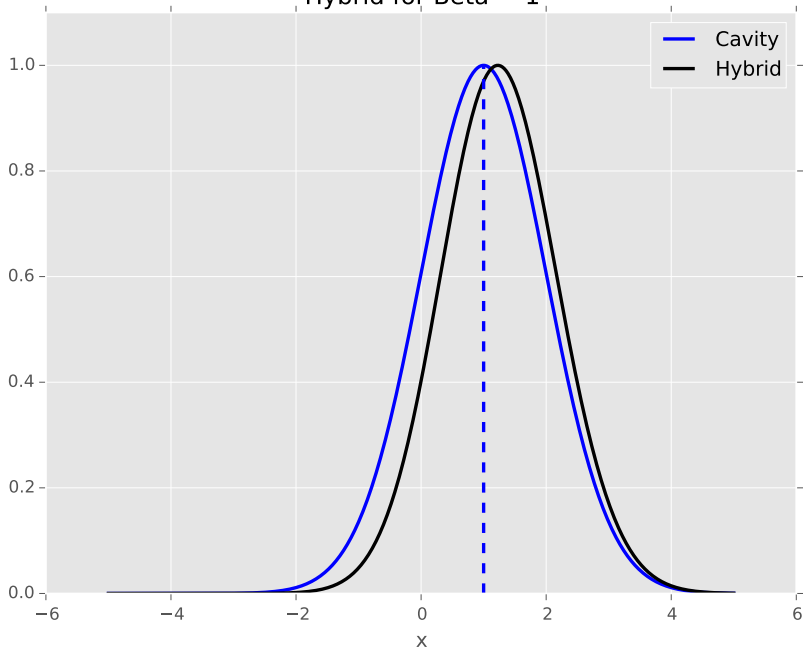
Target factor



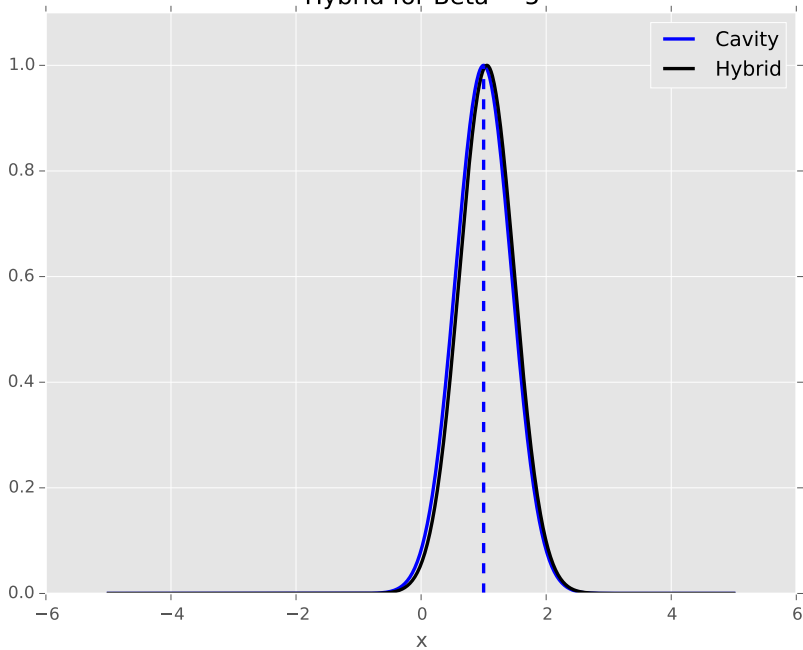
Hybrid for Beta = 0.5



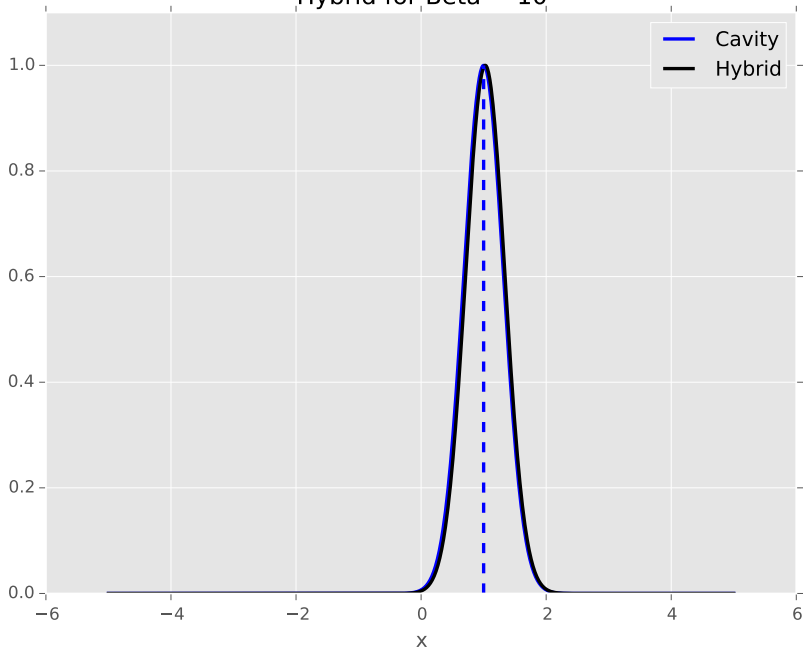
Hybrid for Beta = 1



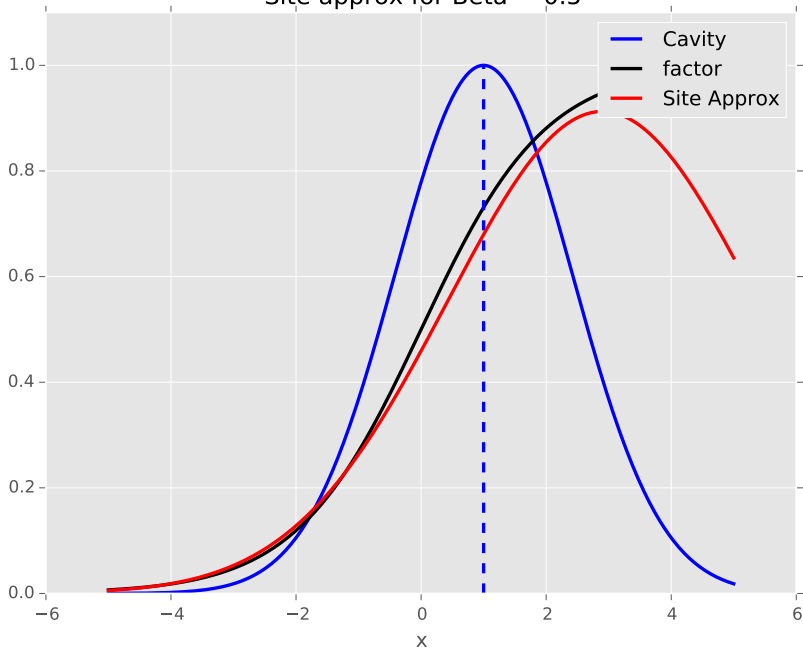
Hybrid for Beta = 5



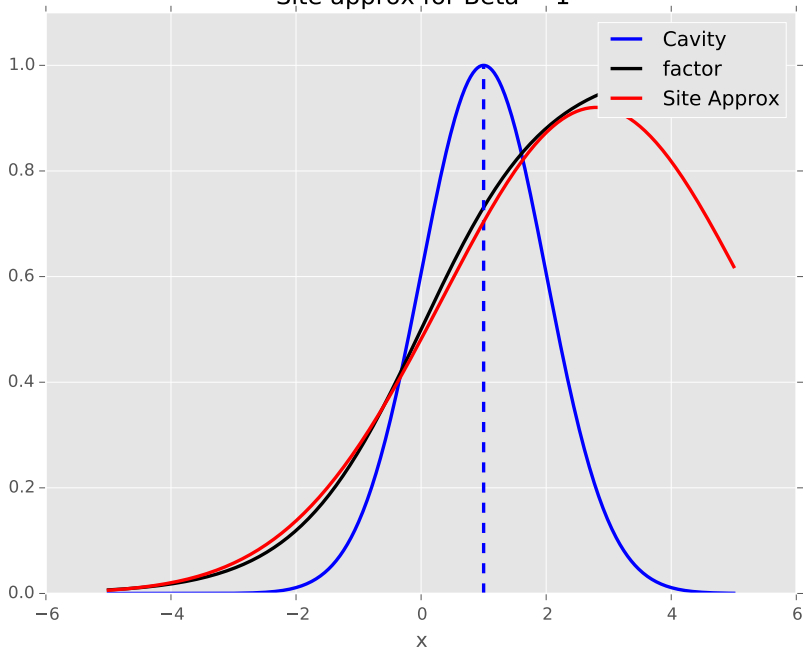
Hybrid for Beta = 10



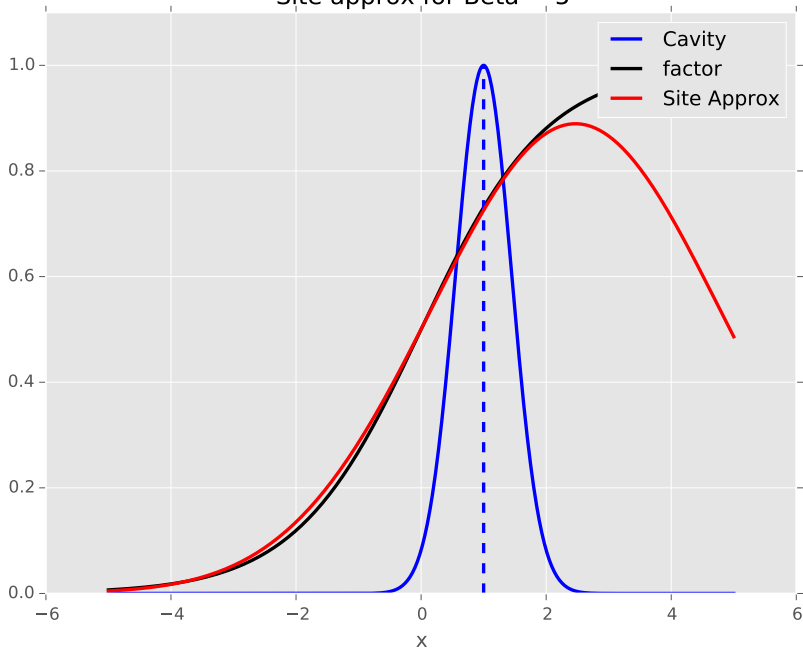
Site approx for Beta = 0.5



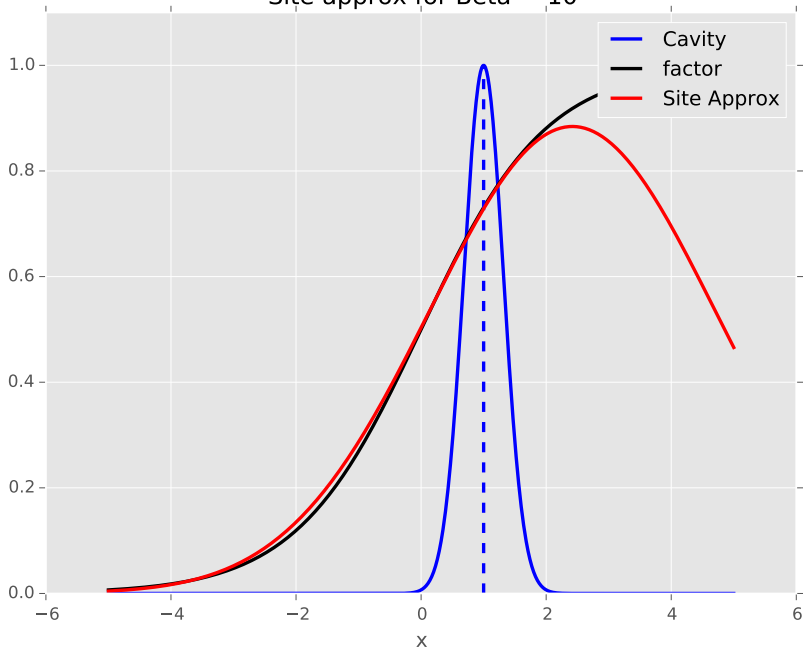
Site approx for Beta = 1



Site approx for Beta = 5



Site approx for Beta = 10



Limit behavior of the approximation

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Theorem (Limit behavior of the factor approximation)

When $\beta_{-i} \rightarrow \infty$, the limit of the EP approximation of $f_i \propto \exp(-\phi_i)$ is:

$$g_i^\infty \propto \exp \left(-\phi'_i(\mu_{-i})(x - \mu_{-i}) - \frac{\phi''(\mu_{-i})}{2} (x - \mu_{-i})^2 \right)$$

- **Many important details in the error term:** non-uniform convergence in μ_{-i}, \dots

Limit behavior of EP iterations

- The limit behavior of parallel EP = the sum of the limit behaviors

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Theorem (Limit behavior of EP)

When $\beta_{-i} \rightarrow \infty$ for all i , the limit of the next EP approximation of $p(x) \propto \exp(-\psi(x))$ is:

$$q_{t+1}^{\infty} \propto \exp \left(-\psi'(\mu_t)(x - \mu_t) - \frac{\psi''(\mu_t)}{2} (x - \mu_t)^2 \right)$$

Limit behavior of EP iterations

- The limit behavior of parallel EP = the sum of the limit behaviors

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Theorem (Limit behavior of EP)

When $\beta_{-i} \rightarrow \infty$ for all i , the limit of the next EP approximation of $p(x) \propto \exp(-\psi(x))$ is:

$$q_{t+1}^{\infty} \propto \exp \left(-\psi'(\mu_t)(x - \mu_t) - \frac{\psi''(\mu_t)}{2} (x - \mu_t)^2 \right)$$

- Did you recognize Newton's algorithm ?

Newton's algorithm

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Objective of Newton's algorithm:
 - find the mode x^* (in order to compute the LA)

$$\psi'(x^*) = 0$$

Newton's algorithm

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

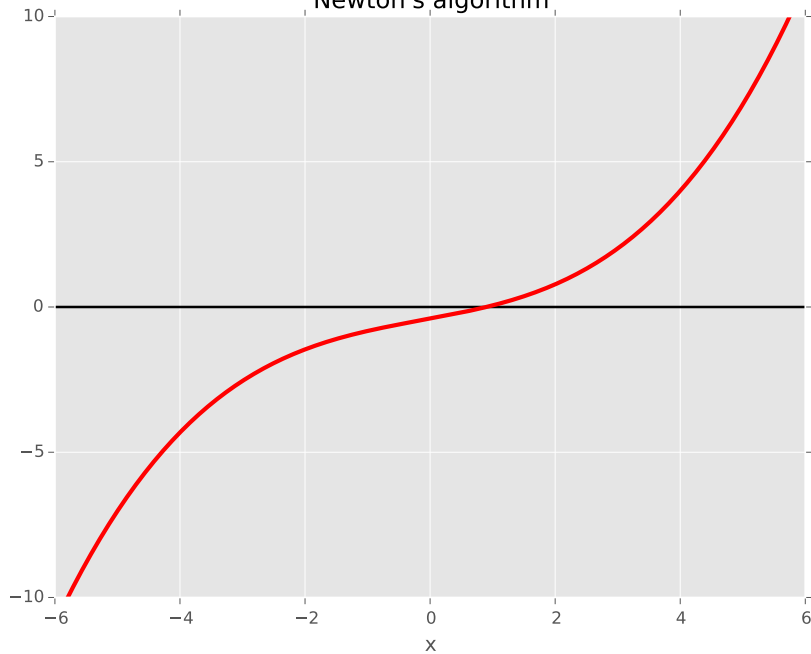
The EP
iteration
behaves like
Newton's
algorithm

- Objective of Newton's algorithm:
 - find the mode x^* (in order to compute the LA)

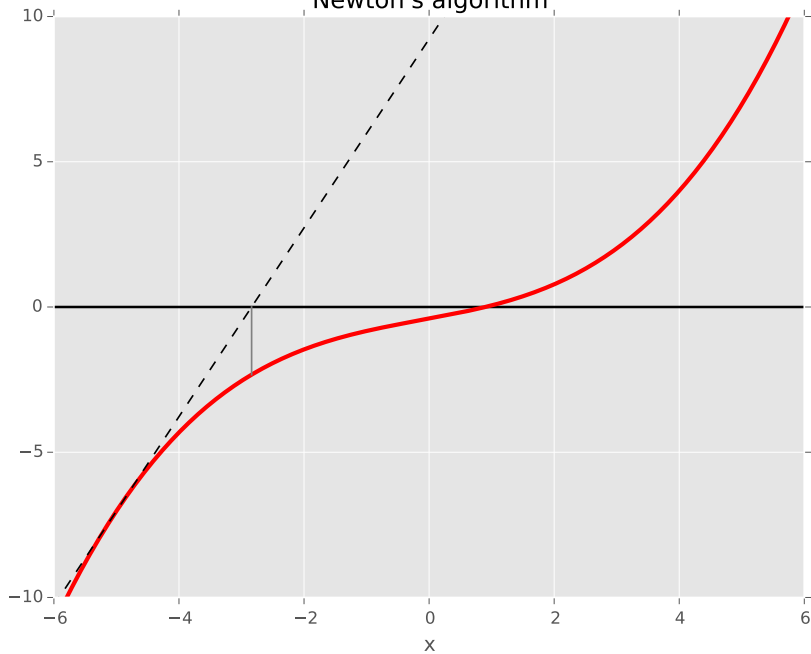
$$\psi'(x^*) = 0$$

- Iterative algorithm:
 - At μ_t , compute the tangent to ψ'
 - Solve tangent = 0: this is μ_{t+1}

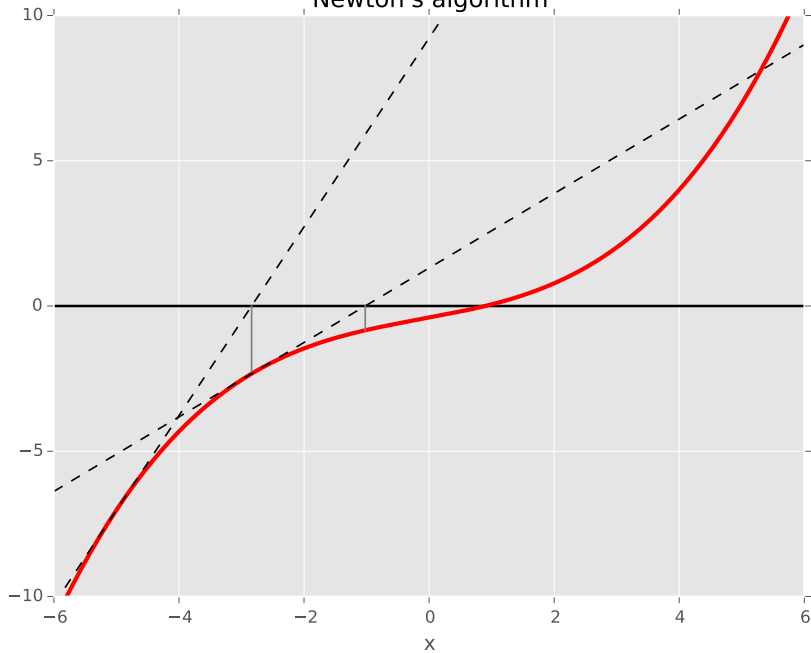
Newton's algorithm



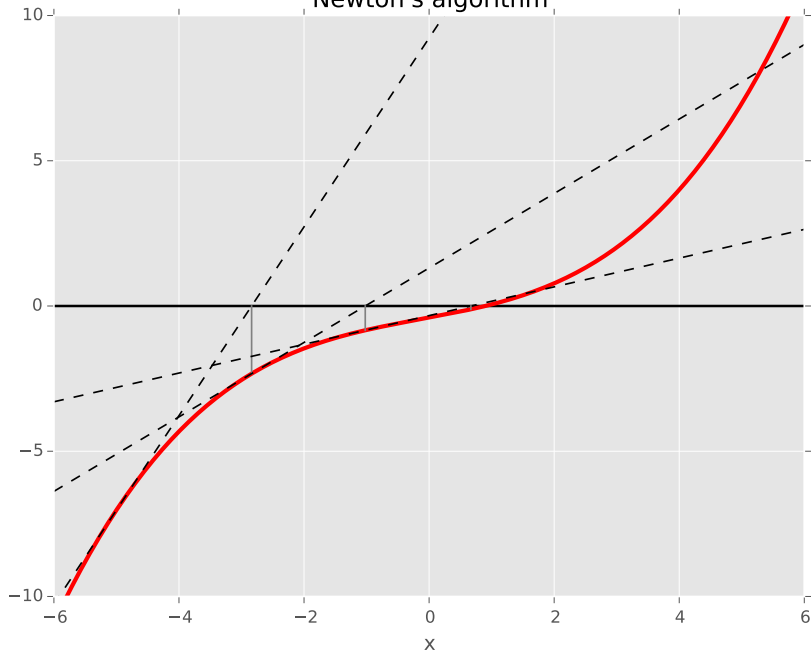
Newton's algorithm



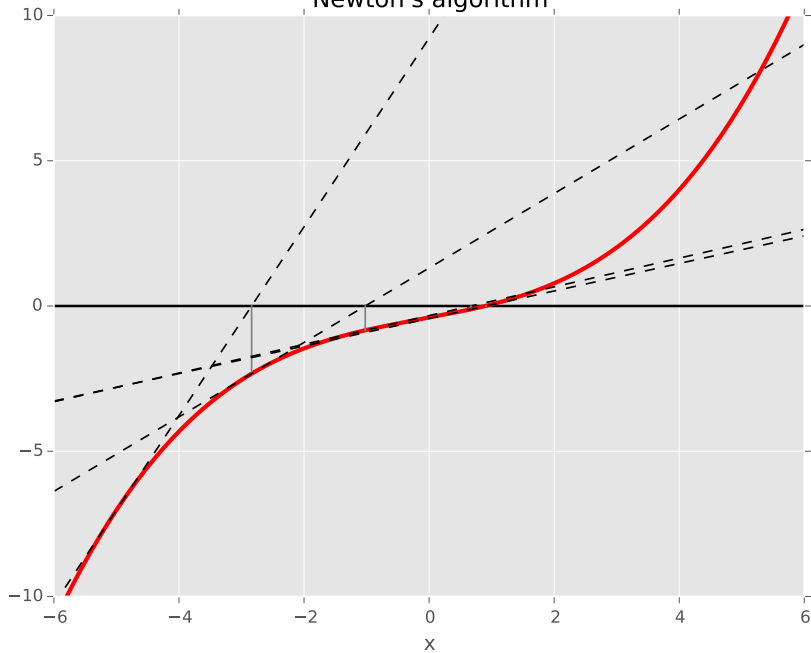
Newton's algorithm



Newton's algorithm



Newton's algorithm



Newton and EP

- At each step, we approximate

$$\psi' \approx ax + b$$

$$\psi \approx a \frac{x^2}{2} + b$$

$$p(x) \approx \exp\left(-a \frac{x^2}{2} - b\right)$$

- Newton is iterating over Gaussian approximations of $p(x)$
!!!

Newton and EP

- At each step, we approximate

$$\psi' \approx ax + b$$

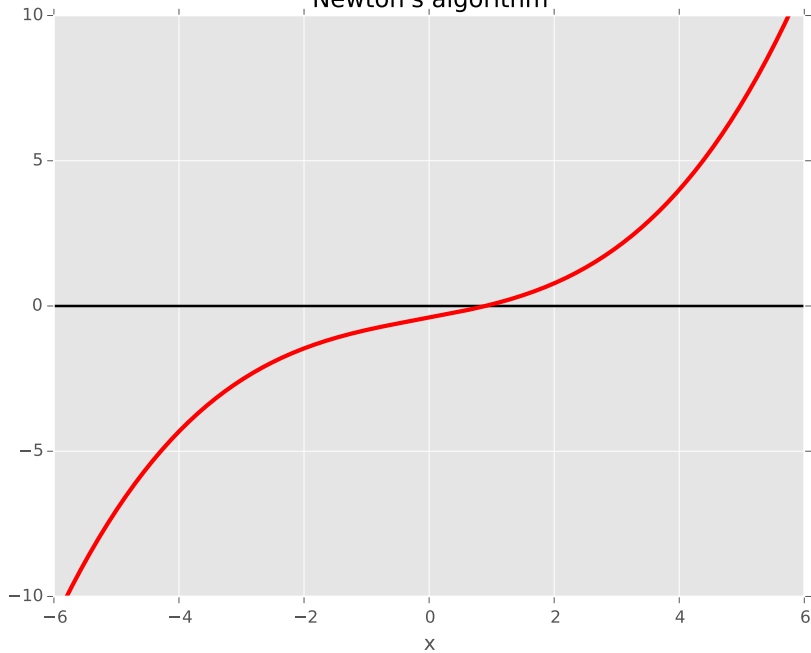
$$\psi \approx a \frac{x^2}{2} + b$$

$$p(x) \approx \exp \left(-a \frac{x^2}{2} - b \right)$$

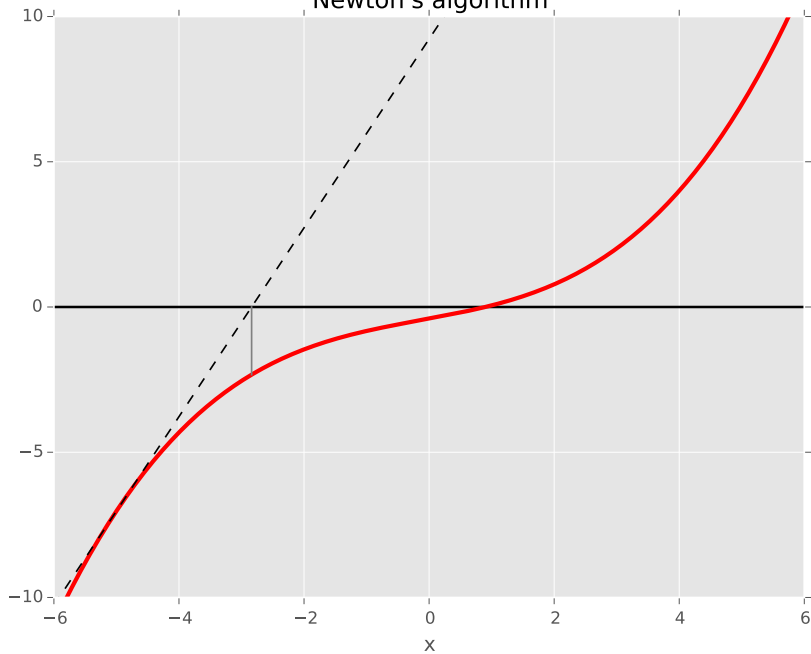
- Newton is iterating over Gaussian approximations of $p(x)$
!!!
- The EP limit approximation is:

$$p(x) \approx \exp \left(-\psi'(\mu_t)(x - \mu_t) - \frac{\psi''(\mu_t)}{2}(x - \mu_t)^2 \right)$$

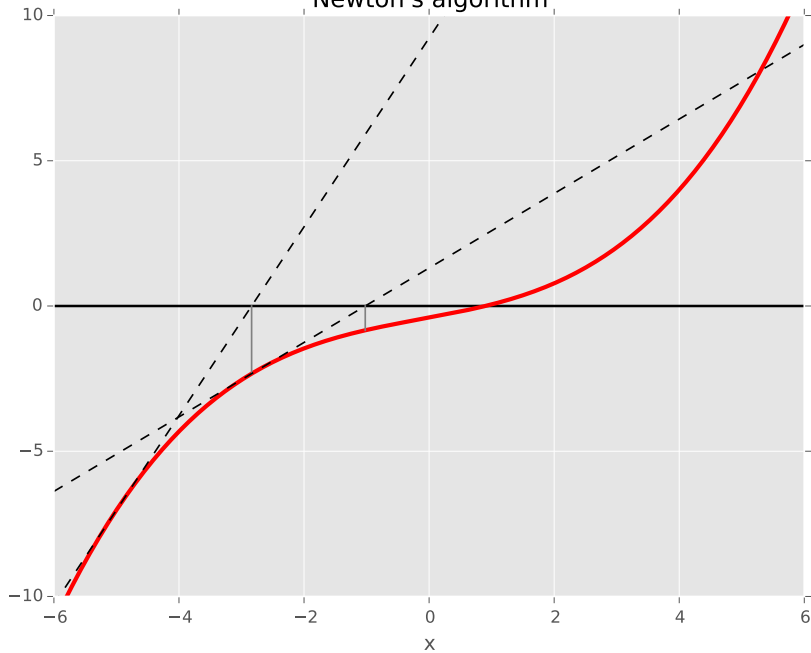
Newton's algorithm



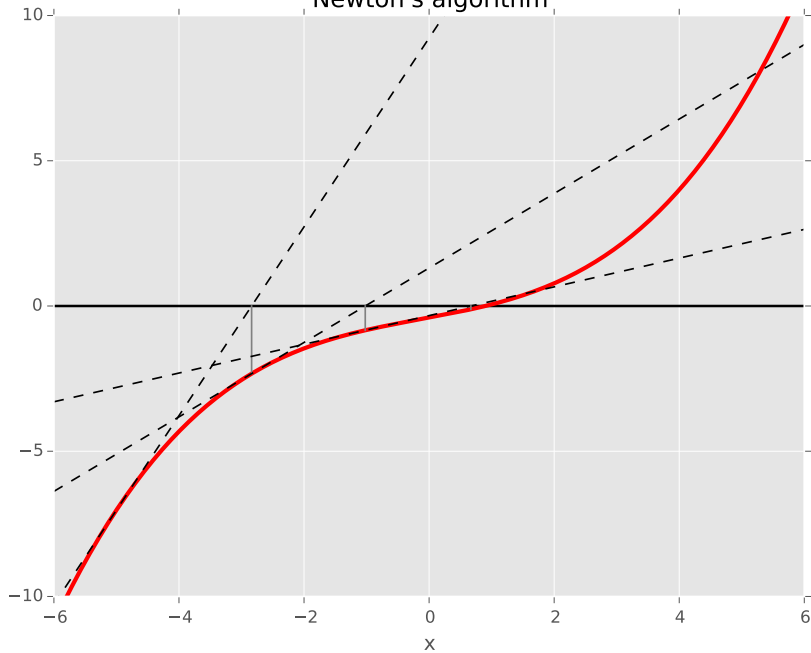
Newton's algorithm



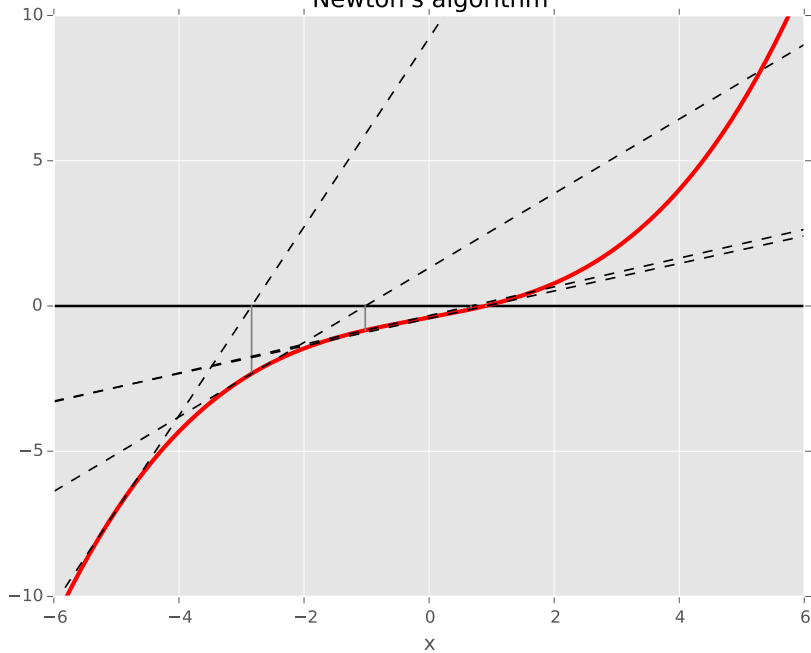
Newton's algorithm



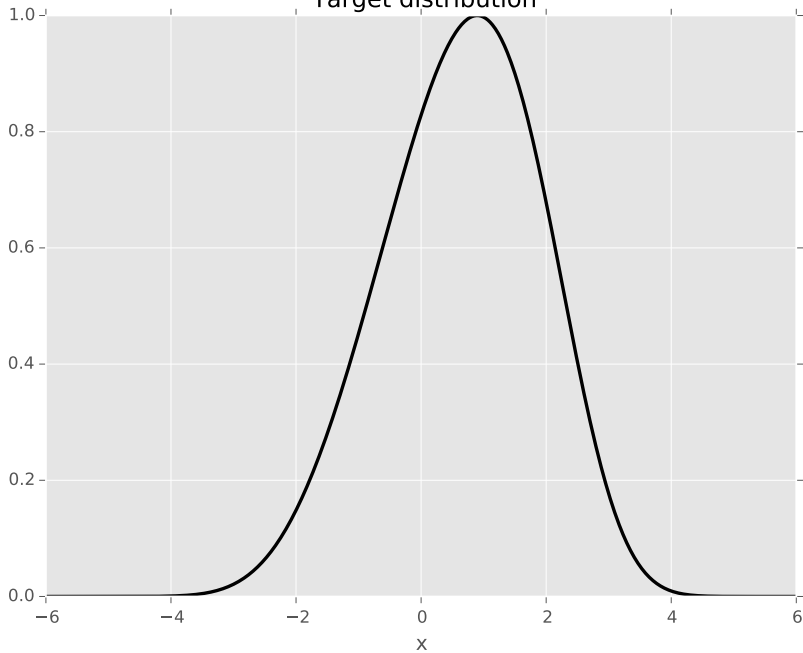
Newton's algorithm



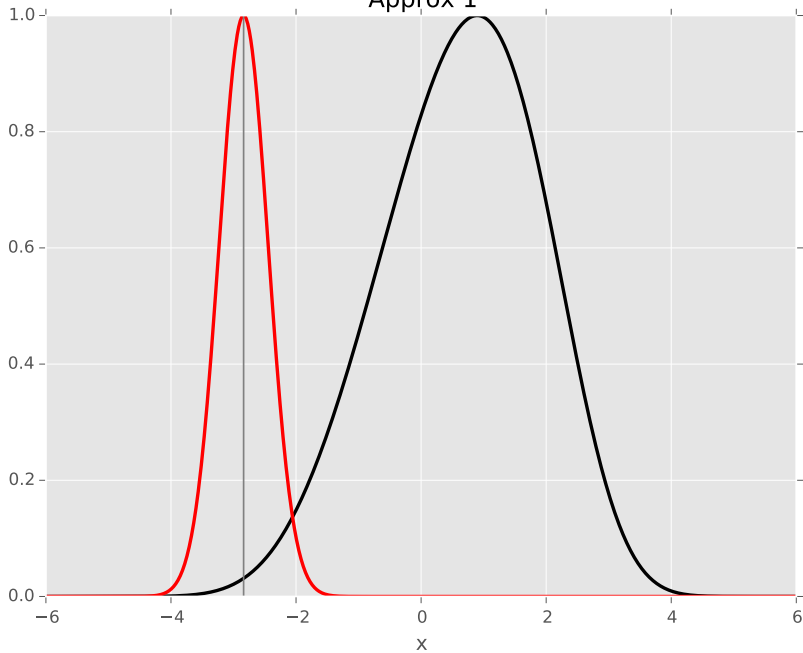
Newton's algorithm



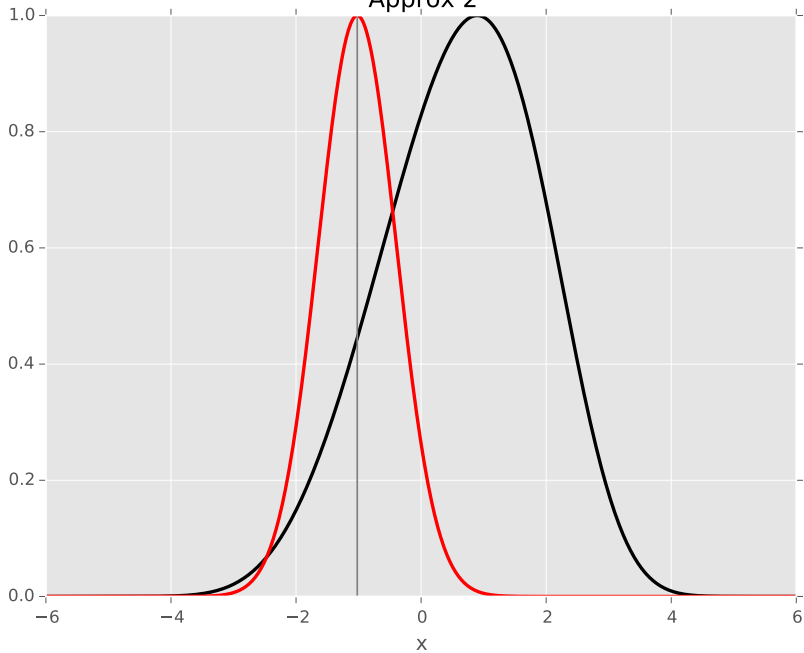
Target distribution



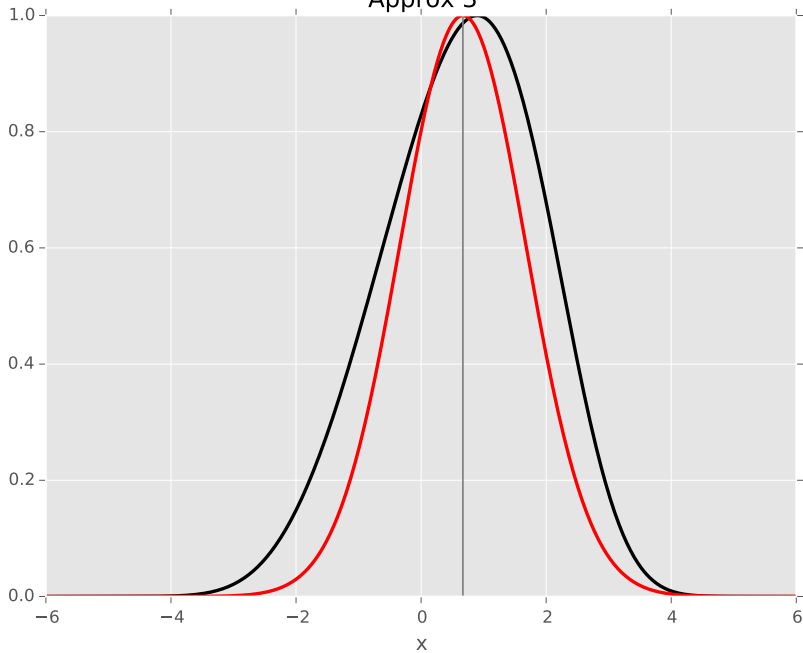
Approx 1



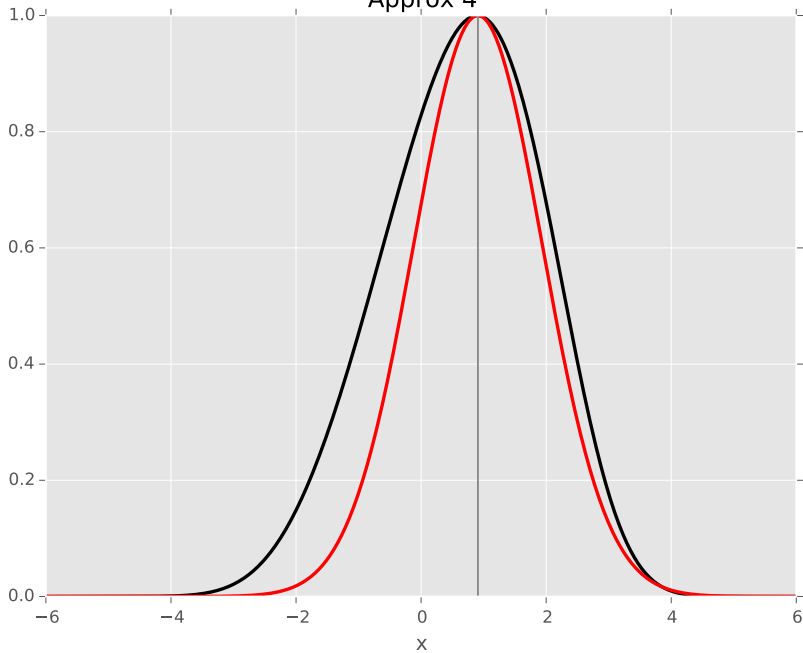
Approx 2



Approx 3



Approx 4



Reaching the high-precision limit

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- High-precision limit \neq large-data limit

Reaching the high-precision limit

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- High-precision limit \neq large-data limit
- Approximation quality result:

$$\text{var}_{p_n} \propto n^{-1}$$

$$v_{EP} \approx \text{var}_{p_n}$$

- Around fixed-points (where it matters), EP is close to Newton

Reaching the high-precision limit

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- High-precision limit \neq large-data limit
- Approximation quality result:

$$\begin{aligned}\text{var}_{p_n} &\propto n^{-1} \\ v_{EP} &\approx \text{var}_{p_n}\end{aligned}$$

- Around fixed-points (where it matters), EP is close to Newton
- We can derive other links
- We can always check

Intuitions from $EP \approx \text{Newton}$

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Intuitively, if EP behaves like Newton in some limit, even away from that limit, it should have similar properties

Intuitions from $EP \approx Newton$

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Intuitively, if EP behaves like Newton in some limit, even away from that limit, it should have similar properties
- Newton is very well-known
 - It converges extremely fast once it gets close to its fixed-point
 - **But** it can fail to converge
 - We have to supplement it with line-search algorithms
 - If we don't, it can “bounce” around its fixed-point

Intuitions from $EP \approx \text{Newton}$

Background

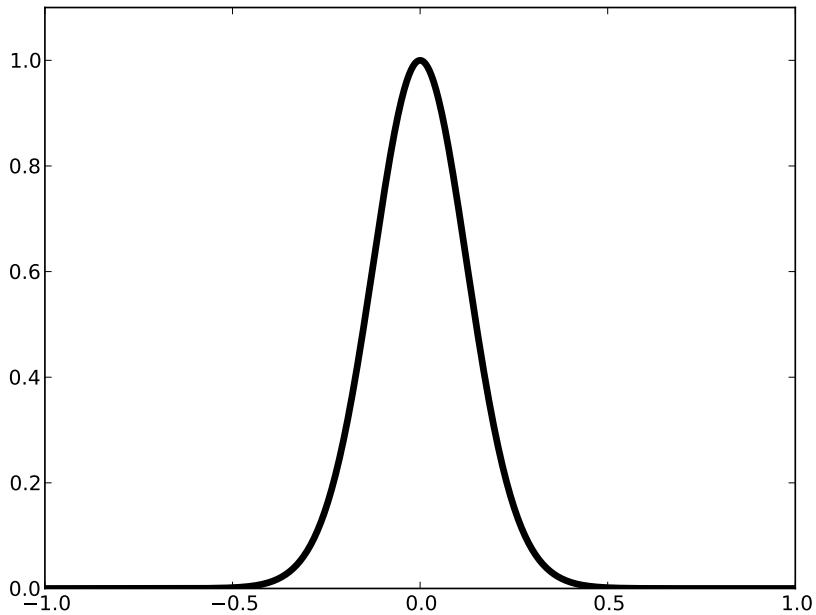
How does EP
work ?

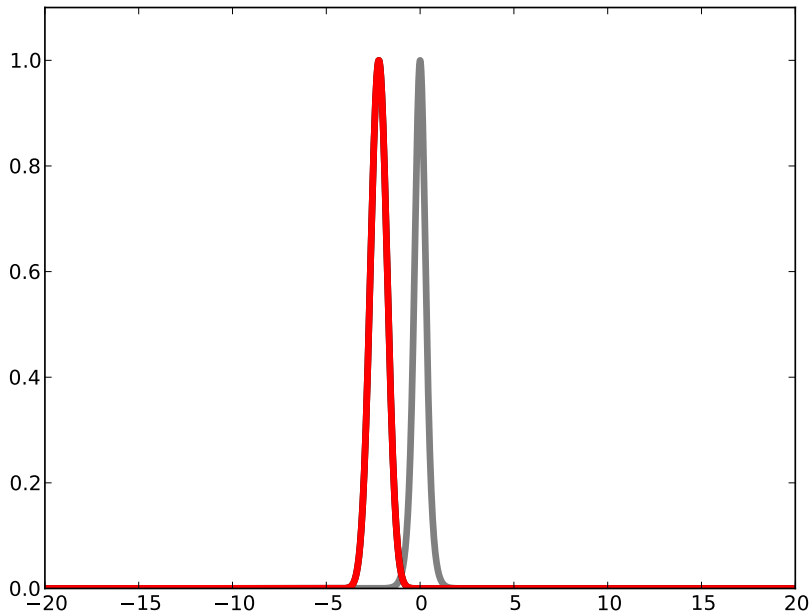
The
large-data
limit

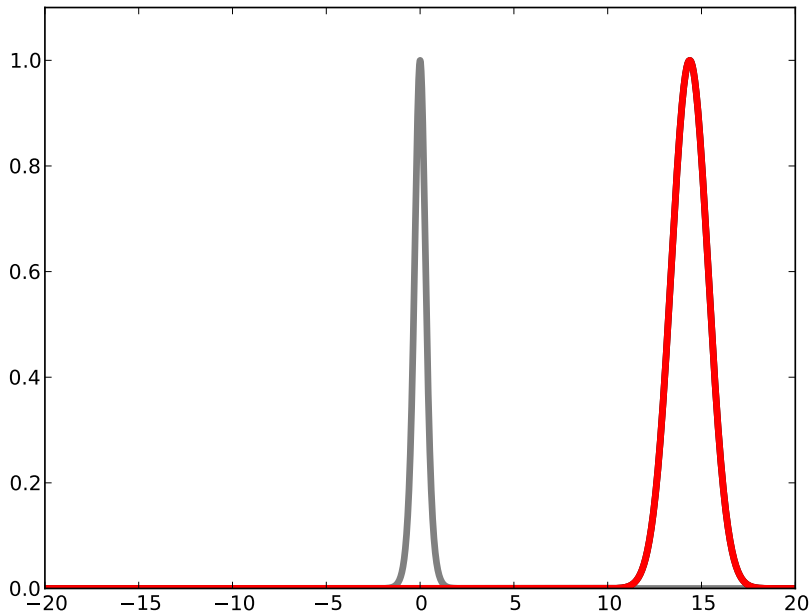
Why does
EP give
accurate ap-
proximations

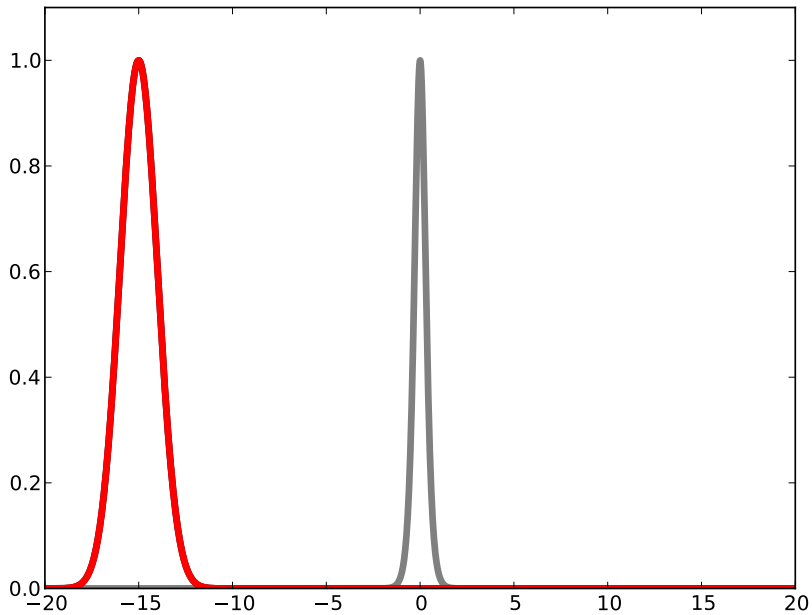
The EP
iteration
behaves like
Newton's
algorithm

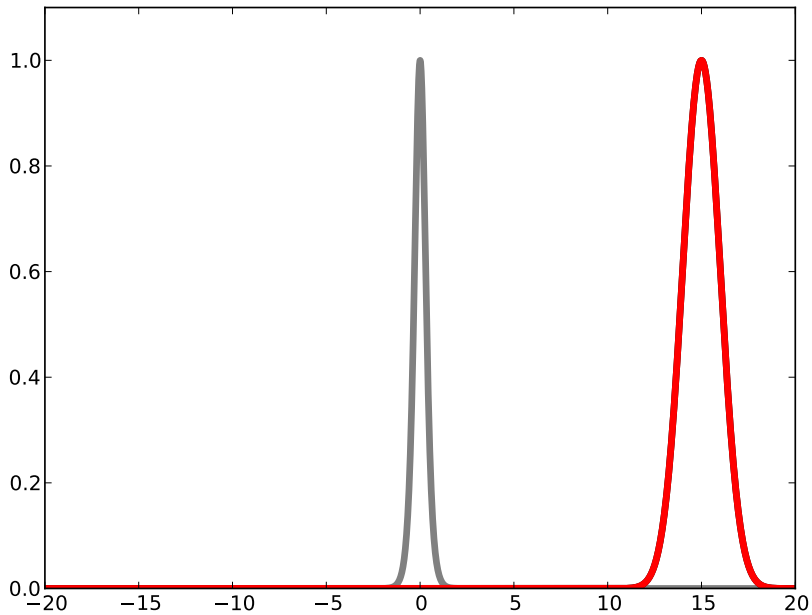
- Intuitively, if EP behaves like Newton in some limit, even away from that limit, it should have similar properties
- Newton is very well-known
 - It converges extremely fast once it gets close to its fixed-point
 - **But** it can fail to converge
 - We have to supplement it with line-search algorithms
 - If we don't, it can “bounce” around its fixed-point
- EP probably has similar properties !

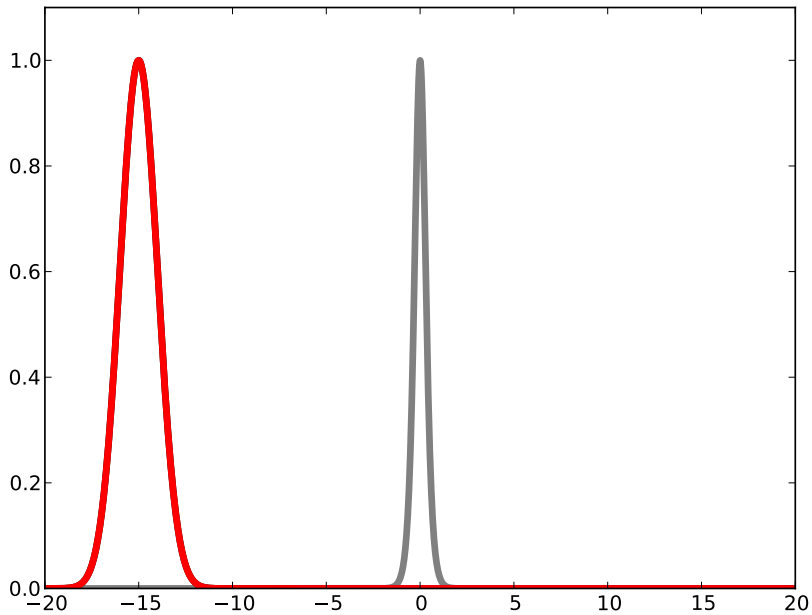


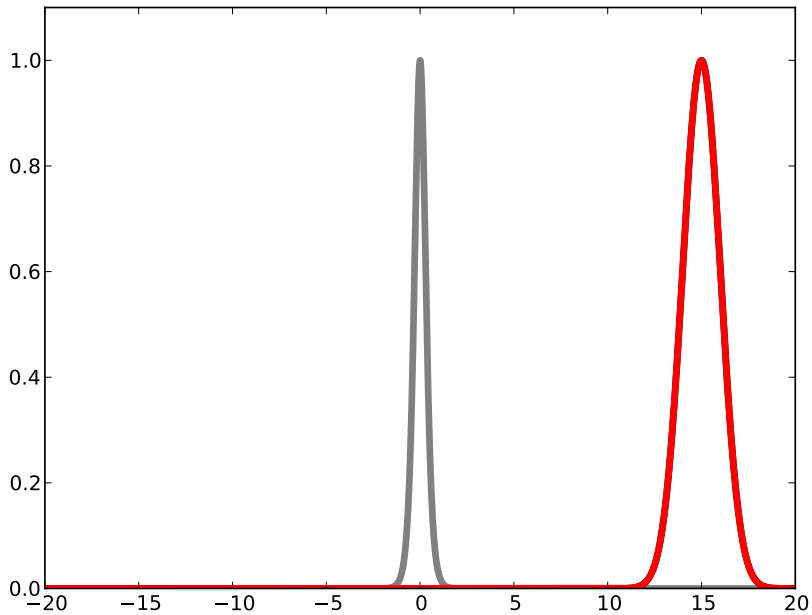


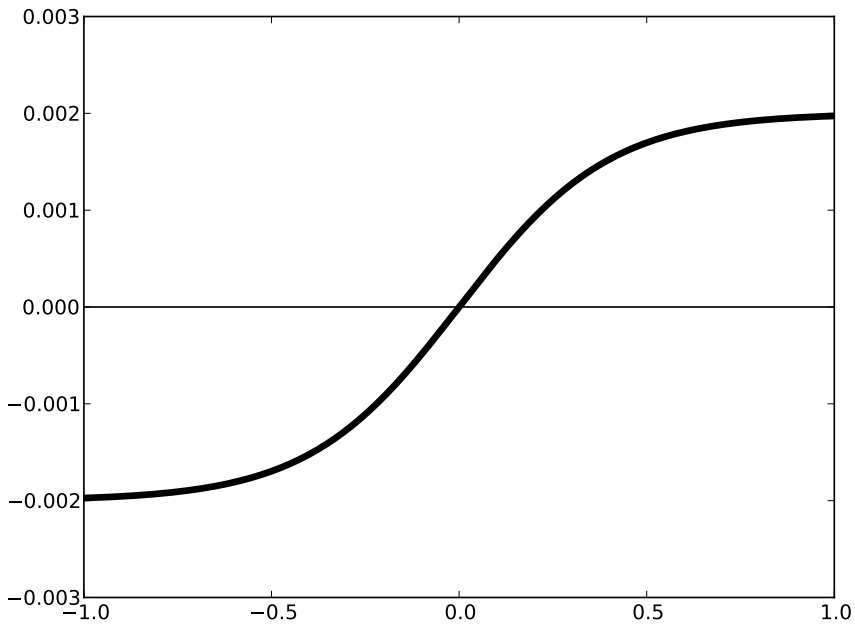












Intuitions from $EP \approx Newton$

Background

How does EP
work ?

The
large-data
limit

- On a multi-modal $p(x)$, Newton has multiple fixed-points

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

Intuitions from $EP \approx Newton$

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- On a multi-modal $p(x)$, Newton has multiple fixed-points
- We can prove that sufficiently peaked modes have an EP fixed-point
 - EP can be “captured” by a mode and miss most of the probability mass
 - Avoid multi-modal distributions like the plague

Summary of our results

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- EP is a better approximation than LA (with some caveats)

Summary of our results

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- EP is a better approximation than LA (with some caveats)
- EP behaves like Newton's algorithm in the high-precision limit
- The high-precision limit is reached in the large-data regime

Summary of our results

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- EP is a better approximation than LA (with some caveats)
- EP behaves like Newton's algorithm in the high-precision limit
- The high-precision limit is reached in the large-data regime
- This sheds new light on the dynamical behavior of EP:
 - It can bounce around it's fixed-point
 - We might need to supplement EP with line-search algorithms
 - EP can be captured by modes

References

Expectation
Propagation
in the large
data limit

Guillaume
Dehaene,
Simon
Barthelmé

Background

How does EP
work ?

The
large-data
limit

Why does
EP give
accurate ap-
proximations

The EP
iteration
behaves like
Newton's
algorithm

- Our work:
 - “Bounding errors of Expectation-Propagation”, Dehaene, Barthelmé, 2015, NIPS
 - “Expectation Propagation is Newton-like in the large-data limit”, Dehaene, Barthelmé, 2015, In review
- Further references:
 - “Birth” of EP: Minka, 2001, UAI
 - Best explanation: Seeger, 2008, Berkely course notes
 - *EP as a way of life*: Gelman, Vehtari, Jylanki, Robert, Chopin, Cunningham