# Probabilistic Machine Learning
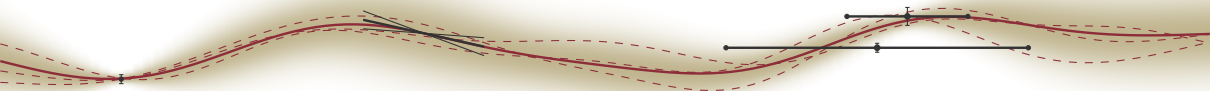## Lecture 24
## Variational Inference

Philipp Hennig

23 July 2023

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING

Setting: Want to find *maximum likelihood* (or MAP) estimate for a model involving a **latent** variable

$$\theta_* = \arg\max_\theta \left[\log p(x \mid \theta)\right] = \arg\max_\theta \left[\log \left(\int p(x, z \mid \theta)\, dz\right)\right]$$

Algorithm: Initialize $\theta_0$, then iterate:

1. Compute $q(z) = p(z \mid x, \theta_{\text{old}})$, thereby setting $D_{\text{KL}}(q \| p(z \mid x, \theta)) = 0$
2. Set $\theta_{\text{new}}$ to the **Maximize** the **Evidence Lower Bound**

$$\theta_{\text{new}} = \arg\max_\theta \mathcal{L}(q, \theta) = \arg\max_\theta \int q(z) \log \left(\frac{p(x, z \mid \theta)}{q(z)}\right)\, dz$$

3. Check for convergence of either the log likelihood, or $\theta$.

UNIVERSITÄT
TÜBINGEN
EBERHARD KARLS

▶ Want to maximize, as function of $\theta := (\pi_k, \mu_k, \Sigma_k)_{k=1,\dots,K}$

$$\log p(\boldsymbol{x} \mid \pi, \mu, \Sigma) = \sum_{n=1}^{N} \log \left( \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n; \mu_k, \Sigma_k) \right)$$

▶ Instead, maximizing the "complete data" likelihood is easier:

$$\log p(\boldsymbol{x}, \boldsymbol{z} \mid \pi, \mu, \Sigma) = \log \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{z_{nk}} \mathcal{N}(x_n; \mu_k, \Sigma_k)^{z_{nk}}$$

$$= \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \underbrace{(\log \pi_k + \log \mathcal{N}(x_n; \mu_k, \Sigma_k)))}_{\text{easy to optimize (exponential families!)}}$$

```python
1  K = 2; (N, D) = X.shape
2
3  pi = np.ones(K) / K # Initialize parameters
4  mu = multivariate_normal.rvs(mean=X.mean(axis=0), cov=np.cov(X, rowvar=False), size=K, random_state=1)
5  sigma = [np.cov(X, rowvar=False) / K for _ in range(K)]
6  r = np.zeros([N, K])
7
8  while True:
9      # E-step
10     for k in range(K):
11         r[:, k] = pi[k] * multivariate_normal.pdf(X, mu[k, :], sigma[k])
12     r /= np.sum(r, axis=1, keepdims=True) + 1e-9
13
14     # M-step
15     Nk = np.sum(r, axis=0)
16     pi = Nk / N
17     mu_old = mu.copy()
18     mu = np.dot(r.T, X) / Nk[:, None]
19     for k in range(K):
20         Xc = X - mu[k]
21         sigma[k] = np.dot(r[:, k] * Xc.T, Xc) / Nk[k]
22
23     if np.linalg.norm(mu - mu_old) < 1e-3:
24         break
```

▶ What if we can not compute the posterior $p(z \mid x, \theta)$ analytically?

▶ In Lecture 22 we saw that EM is iteratively optimizing a lower bound on the log likelihood

$$\log p(x) = \underbrace{\int q(z) \log \left( \frac{p(x, z)}{q(z)} \right) dz}_{\mathcal{L}(q)} - \underbrace{\int q(z) \log \left( \frac{p(z \mid x)}{q(z)} \right)}_{D_{\mathsf{KL}}(q\|p(z|x))}$$

in which we set $D_{\mathsf{KL}}(q\|p(z \mid x)) = 0$ at each iteration.

▶ More generally, we could consider any $q$ and optimize the **Evidence Lower Bound** (ELBO) $\mathcal{L}(q)$

▶ Because $p(x)$ is a constant (assuming fixed model), maximizing the ELBO amounts to minimizing the KL divergence between $q$ and the posterior $p(z \mid x)$, and thus **finding a good approximation to the posterior**.

▶ In principle, this is an **optimization in the space of probability distributions** $q$. It is known as **Variational Inference**.

▶ in general, maximizing $\mathcal{L}(q)$ wrt. $q(z)$ is hard, because the extremum is exactly at $q(z) = p(z \mid x)$

▶ we could choose a **parametric** approximation $q(z \mid \lambda)$, and optimize wrt. $\lambda$. This is restrictive, and does not guarantee success, because we still need to be able to evaluate the integrals

$$\mathcal{L}(q) = \int q(z) \log \left( \frac{p(x, z)}{q(z)} \right) \, dz \qquad \text{or} \qquad D_{\mathsf{KL}}(q(z) \| p(z \mid x)) = \int q(z) \log \left( \frac{q(z)}{p(z \mid x)} \right) \, dz$$

▶ How does one find a *function* (even more, a probability distribution) that minimizes a *functional*?

▶ let's assume that $q(z)$ **factorizes**

$$q(z) = \prod_i q_i(z_i)$$

▶ then the bound simplifies. Let's focus on one particular variable $z_j$:

$$\mathcal{L}(q) = \int \prod_i q_i(z_i) \left( \log p(x, z) - \sum_i \log q_i(z_i) \right) dz$$

$$= \int q_j(z_j) \left( \int \log p(x, z) \prod_{i \neq j} q_i(z_i) \, dz_i \right) dz_j - \int q_j(z_j) \log q_j(z_j) dz_j + \text{const.}$$

$$= \int q_j(z_j) \log \tilde{p}(x, z_j) \, dz_j - \int q_j(z_j) \log q_j(z_j) \, dz_j + \text{const.} \quad \text{where}$$

$$\log \tilde{p}(x, z_j) := \mathbb{E}_{q, i \neq j}[\log p(x, z)] + \text{const.}$$

this is a **function** of $z_j$, and a lower bound on $\log \tilde{p}(x) = \log \int \tilde{p}(x, z_j) \, dz_j$.

Consider a joint distribution $p(x, z) = \tilde{p}(x, z) \cdot \text{const.}$ To find a "good" but tractable approximation $q(z)$,

▶ **assume $q$ factorizes** $q(z) = \prod_i q_i(z_i)$. Initialize all $q_i$ to some initial *distribution*.

▶ Iteratively compute

$$\mathcal{L}(q) = \int q_j \log \tilde{p}(x, z_j) \, dz_j - \int q_j \log q_j \, dz_j + \text{const.} = -D_{\mathsf{KL}}(q_j(z) \| \tilde{p}(x, z_j)) + \text{const.}$$

and maximize wrt. $q_j$. Doing so *minimizes* $D_{\mathsf{KL}}(q(z_j) \| \tilde{p}(x, z_j))$, thus the minimum is at $q_j^*$ with

$$\log q_j^*(z_j) = \log \tilde{p}(x, z_j) = \mathbb{E}_{q, i \neq j}(\log p(x, z)) + \text{const.} \qquad (\star)$$

this expression identifies a **function** $q_j$ of $z_j$, not some parametric form.

▶ the optimization converges, because $-\mathcal{L}(q)$ can be shown to be *convex* wrt. $q$.

In physics, this trick is known as **mean field theory** (because an *n*-body problem is separated into *n* separate problems of individual particles who are affected by the "mean field" $\tilde{p}$ summarizing the expected effect of all other particles).

► is a general framework to construct approximating **probability distributions** $q(z)$ to non-analytic posterior distributions $p(z \mid x)$ by minimizing the **functional**

$$q^* = \underset{q \in \mathcal{Q}}{\arg\min}\, D_{KL}(q(z) \| p(z \mid x)) = \underset{q \in \mathcal{Q}}{\arg\max}\, \mathcal{L}(q)$$

► the beauty is that we get to *choose q*, so one can nearly always find a tractable approximation.

► If we impose the *mean field approximation* $q(z) = \prod_i q(z_i)$, get

$$\log q_j^*(z_j) = \mathbb{E}_{q, i \neq j}(\log p(x, z)) + \text{const.}.$$

► for Exponential Family $p$, we need the expectation under $q$ of the sufficient statistics.

Variational Inference is an extremely flexible and powerful approximation method. Its downside is that constructing the bound and update equations can be tedious. For a quick test, variational inference is often not a good idea. But for a deployed product, it can be the most powerful tool in the box.

▶ **Variational Inference:** is a general framework to construct approximating **probability distributions** $q(z)$ to non-analytic posterior distributions $p(z \mid x)$ by minimizing the **functional**

$$q^* = \underset{q \in \mathcal{Q}}{\arg\min} \, D_{KL}(q(z) \| p(z \mid x)) = \underset{q \in \mathcal{Q}}{\arg\max} \, \mathcal{L}(q)$$

▶ If we impose the *mean field approximation* $q(z) = \prod_i q(z_i)$, the updates reduce to a form of "coordinate ascent in distribution space"

$$\log q_j^*(z_j) = \mathbb{E}_{q, i \neq j}(\log p(x, z)) + \text{const.}$$

▶ in practice, the construction of the bound involves
  1. write down the log joint $\log p(x, z)$, decide on a factorization $q(z) = \prod_i q(z_i)$
  2. inspect the algebraic form of $\log q_j^*(z_j) = \mathbb{E}_{q, i \neq j}(\log p(x, z)) + \text{const.}$ and identify the distribution $q_j$
  3. once all $q_j^*$ are identified by type, try to find analytic expressions for $\mathbb{E}_{q, i \neq j}(\log p(x, z))$.
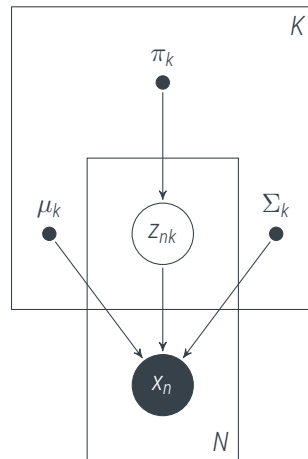
# Example: The Gaussian Mixture Model

► Remember EM for Gaussian mixtures $\theta := (\pi, \mu, \Sigma)$

$$p(x, z \mid \mu, \Sigma, \pi) = \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{z_{nk}} \cdot \mathcal{N}(x_n; \mu_k, \Sigma_k)^{z_{nk}}$$

$$= \prod_{n=1}^{N} p(z_{n:} \mid \pi) \cdot p(x_n \mid z_{n:}, \mu, \Sigma)$$

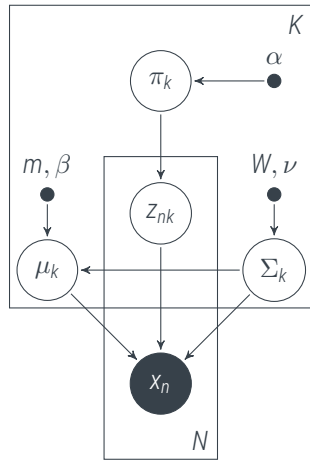▶ Remember EM for Gaussian mixtures $\theta := (\pi, \mu, \Sigma)$

$$p(x, z \mid \mu, \Sigma, \pi) = \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{z_{nk}} \cdot \mathcal{N}(x_n; \mu_k, \Sigma_k)^{z_{nk}}$$

▶ For Bayesian inference, turn parameters into variables

$$p(x, z, \pi, \mu, \Sigma) = p(x, z \mid \pi, \mu, \Sigma) \cdot p(\pi) \cdot p(\mu \mid \Sigma) \cdot p(\Sigma)$$

$$p(\pi) = \mathcal{D}(\pi \mid \alpha) = \frac{\Gamma\left(\sum_{k=1}^{K} \alpha_k\right)}{\prod_k \Gamma(\alpha_k)} \prod_k \pi_k^{\alpha_k - 1}$$

$$p(\mu \mid \Sigma) \cdot p(\Sigma) = \prod_{k=1}^{K} \mathcal{N}(\mu_k; m, \Sigma_k / \beta) \cdot \mathcal{W}(\Sigma_k^{-1}; W, \nu)$$

# Constructing the Variational Approximation

▶ We know that the full posterior $p(z, \pi, \mu, \Sigma \mid x)$ is intractable (check the graph!)
▶ But let's consider an approximation $q(z, \pi, \mu, \Sigma)$ with the factorization

$$q(z, \pi, \mu, \Sigma) = q(z) \cdot q(\pi, \mu, \Sigma)$$

▶ from $(\star)$, we have

$$
\begin{aligned}
\log q^*(z) &= \mathbb{E}_{q(\pi, \mu, \Sigma)} \left( \log p(x, z, \pi, \mu, \Sigma) \right) + \text{const.} \\
&= \mathbb{E}_{q(\pi)} \left( \log p(z \mid \pi) \right) + \mathbb{E}_{q(\mu, \Sigma)} \left( \log p(x \mid z, \mu, \Sigma) \right) + \text{const.} \\
&= \sum_n^N \sum_k^K z_{nk} \underbrace{\left( \mathbb{E}_{q(\pi)}(\log \pi_k) + \frac{1}{2} \mathbb{E}_{q(\mu, \Sigma)}(\log |\Sigma^{-1}| - (x_n - \mu_k)^\mathsf{T} \Sigma_k^{-1}(x - \mu_k)) \right)}_{=: \log \rho_{nk}} + \text{const.}
\end{aligned}
$$

$$q^*(z) \propto \prod_n \prod_k \rho_{nk}^{z_{nk}} \qquad \text{define } r_{nk} = \frac{\rho_{nk}}{\sum_{j=1}^K \rho_{nj}}, \text{ then } q^*(z) = \prod_n \prod_k r_{nk}^{z_{nk}} \text{ with } \mathbb{E}_{q(z)}[z] = r_{nk}$$

# Look through the Math!
Identifying the approximation constructed

UNIVERSITÄT TÜBINGEN

▶ a **discrete** or **multinomial** distribution is defined by

$$q(z) = \prod_{k=1}^{K} \rho_k^{z_k} \quad \text{for } z_k \in \{0, 1\}, \sum_k z_k = 1$$

It has expectations values $\mathbb{E}[z_k] = \rho_k$ and $\mathbb{E}[z_k z_j] = \rho_k \delta_{kj}$

▶ Define some convenient notation:

$$N_k := \sum_{n=1}^{N} r_{nk} \qquad \bar{x}_k := \frac{1}{N_k} \sum_{n=1}^{N} r_{nk} x_n \qquad S_k := \frac{1}{N_k} \sum_{n=1}^{N} r_{nk} (x_n - \bar{x}_k)(x_n - \bar{x}_k)^{\mathsf{T}}$$

▶ from $(\star)$, we have

$$\log q^*(\pi, \mu, \Sigma) = \mathbb{E}_{q(z)} \left( \log p(x, z, \pi, \mu, \Sigma) \right) + \text{const.}$$

$$= \mathbb{E}_{q(z)} \left( \log p(\pi) + \sum_k \log p(\mu_k, \Sigma_k) + \log p(z \mid \pi) + \sum_n \log p(x_n \mid z, \mu, \Sigma) \right)$$

$$= \log p(\pi) + \sum_{k=1}^{K} \log p(\mu_k, \Sigma_k) + \mathbb{E}_{q(z)}(\log p(z \mid \pi))$$

$$+ \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbb{E}(z_{nk}) \log \mathcal{N}(x_n; \mu_k, \Sigma_k) + \text{const.}$$

# Constructing the Variational Approximation

using $q(z, \pi, \mu, \Sigma) = q(z) \cdot q(\pi, \mu, \Sigma)$

$$\log q^*(\pi, \mu, \Sigma) = \log p(\pi) + \sum_{k=1}^{K} \log p(\mu_k, \Sigma_k) + \mathbb{E}_{q(z)}(\log p(z \mid \pi))$$

$$+ \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbb{E}(z_{nk}) \log \mathcal{N}(x_n; \mu_k, \Sigma_k) + \text{const.}$$

▶ The bound exposes an **induced factorization** into $q(\pi, \mu, \Sigma) = q(\pi) \cdot \prod_{k=1}^{K} q(\mu_k, \Sigma_k)$

$$\text{where } \log q(\pi) = \log p(\pi) + \mathbb{E}_{q(z)}(\log p(z \mid \pi)) + \text{const.}$$

$$= (\alpha - 1) \sum_k \log \pi_k + \sum_k \sum_n r_{nk} \log \pi_k + \text{const.}$$

$$q(\pi) = \mathcal{D}(\pi, \alpha_k := \alpha + N_k) \quad \text{with } N_k = \sum_n r_{nk}$$

## Look through the Math!
Identifying the approximation constructed

UNIVERSITÄT
TÜBINGEN

► a **Dirichlet** distribution is defined by

$$p(\pi \mid \alpha) = \mathcal{D}(\pi; \alpha) = \frac{\Gamma(\hat{\alpha})}{\prod_{k=1}^{K} \Gamma(\alpha_k)} \prod_k \pi^{\alpha_k - 1} = \frac{1}{B(\alpha)} \prod_k \pi^{\alpha_k - 1} \qquad \hat{\alpha} := \sum_d \alpha_k$$

It satisfies the following identities:

► $\mathbb{E}_p(\pi_k) = \frac{\alpha_k}{\hat{\alpha}}$

► $\mathrm{var}_p(\pi_k) = \frac{\alpha_k(\hat{\alpha} - \alpha_k)}{\hat{\alpha}^2(\hat{\alpha} + 1)}$

► $\mathrm{cov}(\pi_k, \pi_i) = -\frac{\alpha_k \alpha_i}{\hat{\alpha}^2(\hat{\alpha} + 1)}$

► $\mathrm{mode}(\pi_k) = \frac{\alpha_k - 1}{\hat{\alpha} - D}$

► $\mathbb{E}_p(\log \pi_k) = F(\alpha_k) - F(\hat{\alpha})$

► $\mathbb{H}(p) = -\int p(\pi) \log p(\pi) \, d\pi = -\sum_k (\alpha_k - 1)(F(\alpha_k) - F(\hat{\alpha})) + \log B(\alpha)$

Where $F(z) = \frac{d}{dz} \log \Gamma(z)$ (the "digamma-function").

`scipy.special.digamma(z)`     https://dlmf.nist.gov/5

# Constructing the Variational Approximation

$$\log q^*(\pi, \mu, \Sigma) = \log p(\pi) + \sum_{k=1}^{K} \log p(\mu_k, \Sigma_k) + \mathbb{E}_{q(z)}(\log p(z \mid \pi))$$

$$+ \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbb{E}(z_{nk}) \log \mathcal{N}(x_n; \mu_k, \Sigma_k) + \text{const.}$$

▶ The bound exposes an **induced factorization** into $q(\pi, \mu, \Sigma) = q(\pi) \cdot \prod_{k=1}^{K} q(\mu_k, \Sigma_k)$

where (leaving out some tedious algebra) $q^*(\mu_k, \Sigma_k) = \mathcal{N}(\mu_k; m_k, \Sigma_k/\beta_k) \mathcal{W}(\Sigma_k^{-1}; W_k, \nu_k)$

with $\beta_k := \beta + N_k \qquad m_k := \dfrac{1}{\beta_k}(\beta m + N_k \bar{x}_k) \qquad \nu_k := \nu + N_k$

$$W_k^{-1} := W^{-1} + N_k S_k + \frac{\beta N_k}{\beta + N_k}(\bar{x}_k - m)(\bar{x}_k - m)^{\mathsf{T}}$$

# Constructing the Variational Approximation

closing the loop

▶ Recall from above:

$$\log q^*(z) = \sum_n^N \sum_k^K z_{nk} \underbrace{\left( \mathbb{E}_{q(\pi)}(\log \pi_k) + \frac{1}{2}\mathbb{E}_{q(\mu,\Sigma)}(\log|\Sigma^{-1}| - (x_n - \mu_k)^\mathsf{T}\Sigma_k^{-1}(x - \mu_k)) \right)}_{=:\log \rho_{nk}} + \text{const.}$$

▶ now we can evaluate $\rho_{nk}$, using tabulated identities

$$\log \tilde{\pi}_k := \mathbb{E}_{\mathcal{D}(\pi;\alpha_k)}(\log \pi_k) = F(\alpha_k) - F\left(\sum_k \alpha_k\right)$$

and for the Wishart:

$$\log \tilde{|\Sigma^{-1}|}_k := \mathbb{E}_{\mathcal{W}(\Sigma_k^{-1};W_k,\nu_k)}(\log|\Sigma_k^{-1}|) = \sum_{d=1}^D F\left(\nu_k + \frac{1-d}{2}\right) + D\log 2 + \log|W_k|$$

$$\mathbb{E}_{\mathcal{N}(\mu_k;m_k,\Sigma_k/\beta_k)\mathcal{W}(\Sigma_k^{-1};W_k,\nu_k)}\left((x_n - \mu_k)^\mathsf{T}\Sigma^{-1}(x_n - \mu_k)\right) = D\beta_k^{-1} + \nu_k(x_n - m_k)^\mathsf{T}W_k(x_n - m_k)$$

▶ this yields the update equation

$$\mathbb{E}_q(z_{nk}) = r_{nk} \propto \tilde{\pi}_k |\tilde{\Sigma^{-1}}|^{1/2} \exp \left( -\frac{D}{2\beta_k} - \frac{\nu_k}{2}(x_n - m_k)^\intercal W_k(x_n - m_k) \right)$$
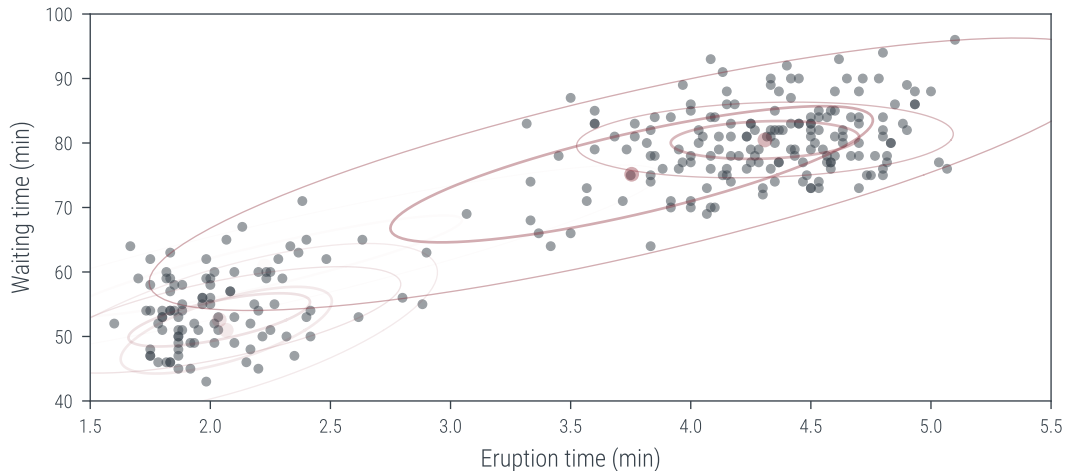
compare this with the EM-update

$$r_{nk} \propto \pi_k |\Sigma^{-1}|^{1/2} \exp \left( -\frac{1}{2}(x_n - \mu_k)^\intercal \Sigma_k^{-1}(x_n - \mu_k) \right)$$

▶ Here, variational Inference is the Bayesian version of EM: Instead of maximizing the likelihood for $\theta = (\mu, \Sigma, \pi)$, we maximize a variational bound.

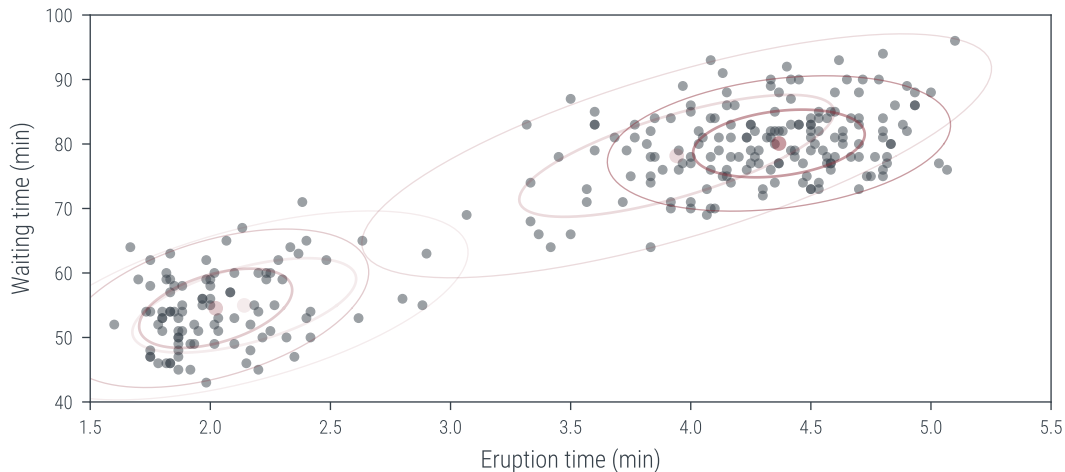▶ One advantage of this is that the posterior can actually "decide" to ignore components:
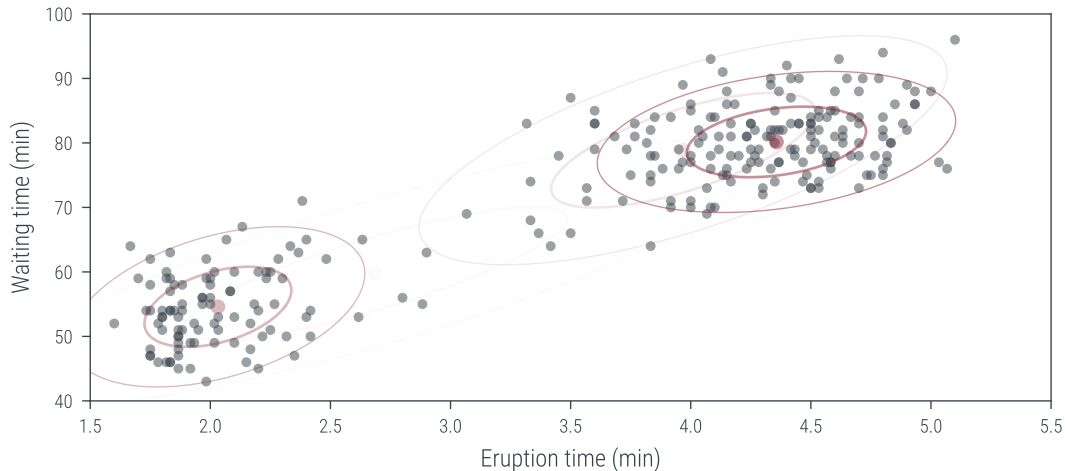
initialized with $K = 6$ components

# Variational Gaussian Mixture Models

initialized with $K = 6$ components

# Variational Gaussian Mixture Models
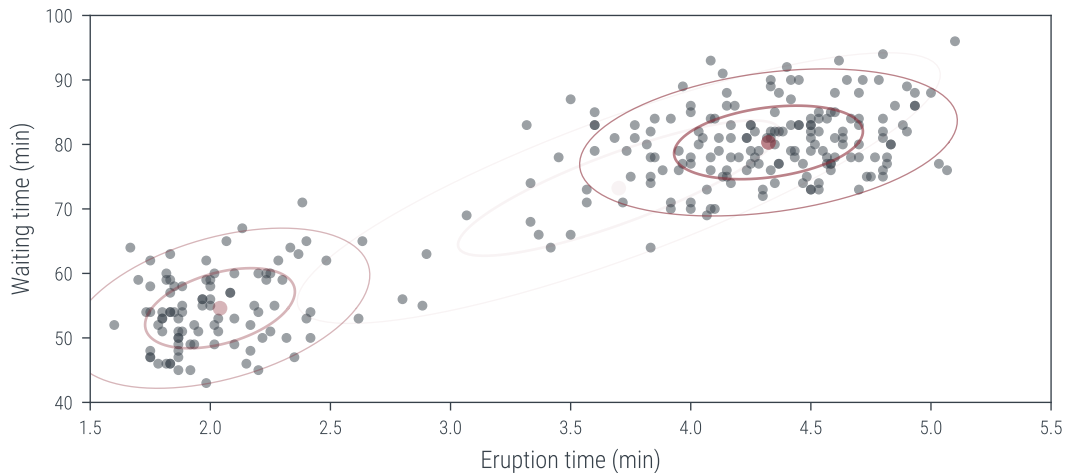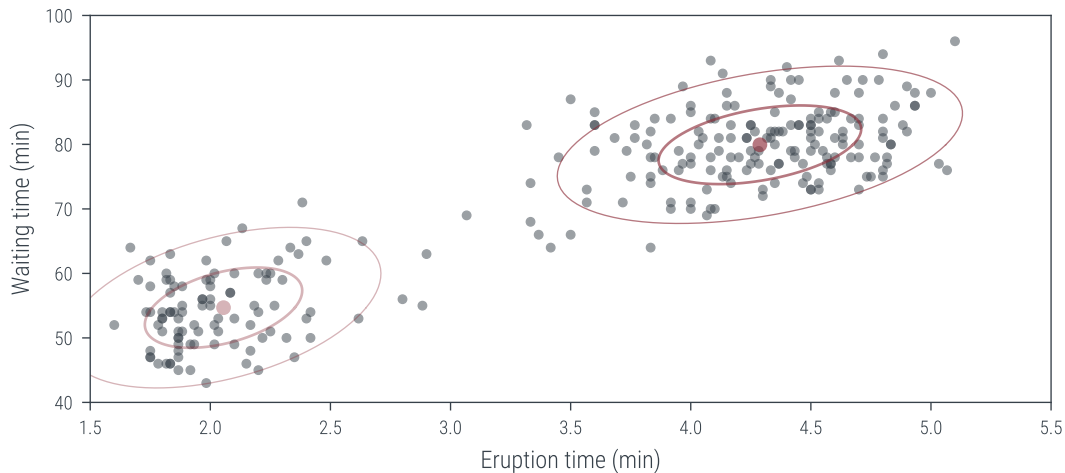
initialized with $K = 6$ components

▶ **Variational Inference:** is a general framework to construct approximating **probability distributions** $q(z)$ to non-analytic posterior distributions $p(z \mid x)$ by minimizing the **functional**

$$q^* = \underset{q \in \mathcal{Q}}{\arg\min} \, D_{KL}(q(z) \| p(z \mid x)) = \underset{q \in \mathcal{Q}}{\arg\max} \, \mathcal{L}(q)$$

▶ If we impose the *mean field approximation* $q(z) = \prod_i q(z_i)$, the updates reduce to a form of "coordinate ascent in distribution space"

$$\log q_j^*(z_j) = \mathbb{E}_{q,i \neq j}(\log p(x,z)) + \text{const.}$$

▶ in practice, the construction of the bound involves
   1. write down the log joint $\log p(x,z)$, decide on a factorization $q(z) = \prod_i q(z_i)$
   2. inspect the algebraic form of $\log q_j^*(z_j) = \mathbb{E}_{q,i \neq j}(\log p(x,z)) + \text{const.}$ and identify the distribution $q_j$
   3. once all $q_j^*$ are identified by type, try to find analytic expressions for $\mathbb{E}_{q,i \neq j}(\log p(x,z))$.

Variational Inference is a powerful mathematical tool to construct efficient approximations to intractable *probability distributions* (not just point estimates, but entire distributions). Often, just imposing factorization is enough to make things tractable. The downside of variational inference is that constructing the bound can take significant ELBOw grease. However, the resulting algorithms are often highly efficient compared to tools that require less derivation work, like Monte Carlo.

"Derive the variational bound while you wait for MCMC to converge."

Please cite this course, as

```
@techreport{Tuebingen_ProbML23,
    title =
        {Probabilistic Machine Learning},
    author = {Hennig, Philipp},
    series = {Lecture Notes
        in Machine Learning},
    year = {2023},
    institution = {Tübingen AI Center}}
```